

Logarithmic barriers for sparse matrix cones

Martin S. Andersen*

Joachim Dahl†

Lieven Vandenberghe*

Abstract

Algorithms are presented for evaluating gradients and Hessians of logarithmic barrier functions for two types of convex cones: the cone of positive semidefinite matrices with a given sparsity pattern, and its dual cone, the cone of sparse matrices with the same pattern that have a positive semidefinite completion. Efficient large-scale algorithms for evaluating these barriers and their derivatives are important in interior-point methods for nonsymmetric conic formulations of sparse semidefinite programs. The algorithms are based on the multifrontal method for sparse Cholesky factorization.

1 Introduction

1.1 Log-det barrier for sparse matrices

We discuss algorithms for evaluating the gradient and Hessian of the ‘log-det’ barrier

$$f(X) = -\log \det X$$

when X is large, sparse, and positive definite. We take f as a function from \mathbf{S}_V^n to \mathbf{R} , where V is the filled sparsity pattern of the Cholesky factor of X , and \mathbf{S}_V^n denotes the set of symmetric matrices with sparsity pattern V . With this convention, and for the standard trace inner product $A \bullet B = \text{tr}(AB)$ of symmetric matrices, the gradient of f at X is

$$\nabla f(X) = -\mathcal{P}(X^{-1}) \tag{1}$$

where \mathcal{P} denotes projection on V , *i.e.*, $\mathcal{P}(A)_{ij} = A_{ij}$ if the pattern V has a nonzero in position i, j and $\mathcal{P}(A)_{ij} = 0$ otherwise. The algorithms presented in this paper exploit properties of filled sparsity patterns (which are also known as *chordal* or *triangulated* patterns) to compute the gradient directly from the Cholesky factor of X , without calculating the rest of the inverse.

The Hessian of f , interpreted as a function from \mathbf{S}_V^n to \mathbf{R} , is defined by

$$\nabla^2 f(X)[Y] = \left. \frac{d}{dt} \nabla f(X + tY) \right|_{t=0} = \mathcal{P}(X^{-1} Y X^{-1}) \quad \forall Y \in \mathbf{S}_V^n. \tag{2}$$

We are interested in efficient methods for evaluating this expression, possibly for multiple matrices Y simultaneously, without computing the entire inverse of X or the products $X^{-1} Y X^{-1}$. We also discuss methods for evaluating the inverse Hessian $\nabla^2 f(X)^{-1}[Y]$.

*Electrical Engineering Department, University of California, Los Angeles. Email: martin.andersen@ucla.edu, vandenbe@ee.ucla.edu.

†MOSEK ApS, Fruebjergvej 3, 2100 København Ø. Email: dahl.joachim@gmail.com.

The function f has an important role as a logarithmic barrier function for the convex cone

$$\mathbf{S}_{V,+}^n = \{X \in \mathbf{S}_V^n \mid X \succeq 0\},$$

the cone of positive semidefinite matrices with sparsity pattern V .

1.2 Conjugate barrier

Equally important is the corresponding dual barrier function

$$f_*(S) = \sup_{X \in \mathbf{S}_V^n} (-S \bullet X - f(X)). \quad (3)$$

This is the conjugate or Legendre transform of f applied to $-S$. The function f_* is a logarithmic barrier function for the dual cone of $\mathbf{S}_{V,+}^n$, which contains the symmetric matrices with sparsity pattern V that have a positive semidefinite completion:

$$\mathbf{S}_{V,c}^n = \mathcal{P}(\mathbf{S}_+^n) = \{\mathcal{P}(X) \mid X \succeq 0\}.$$

The dual barrier $f_*(S)$ can be computed as $f_*(S) = \log \det \hat{X} - n$ where \hat{X} is the maximizer in the definition of f_* , *i.e.*, the solution of the nonlinear equation $\nabla f(X) = -S$ or

$$\mathcal{P}(X^{-1}) = S \quad (4)$$

with variable $X \in \mathbf{S}_V^n$. The gradient and Hessian of f_* at S also follow from the maximizer \hat{X} by applying standard properties of Legendre transforms:

$$\nabla f_*(S) = -\hat{X}, \quad \nabla^2 f_*(S) = \nabla^2 f(\hat{X})^{-1}. \quad (5)$$

For general sparsity patterns V , the maximizer \hat{X} in (3) needs to be computed by iterative methods. For filled patterns V , however, efficient direct algorithms exist. The algorithms discussed in this paper compute a Cholesky factorization of \hat{X} , given the matrix S , using a finite recursion that is very similar and comparable in cost to a Cholesky factorization.

There is an interesting connection between the dual barrier f_* and the maximum determinant positive definite completion problem, which has been extensively studied in linear algebra [GJSW84, Lau01]. The optimization problem in (3) is the Lagrange dual of the convex optimization problem

$$\begin{aligned} & \text{minimize} && -\log \det Z - n \\ & \text{subject to} && \mathcal{P}(Z) = S, \end{aligned} \quad (6)$$

with variable $Z \in \mathbf{S}^n$. The primal and dual optimal solutions \hat{X} and Z are related by the optimality condition $Z^{-1} = \hat{X}$. The solution \hat{X} of (4) is therefore the inverse of the maximum determinant positive definite completion of S .

1.3 Applications

Efficient gradient and Hessian evaluations for f and f_* are critical to the performance of interior-point methods for conic optimization problems associated with the cones $\mathbf{S}_{V,+}^n$ and $\mathbf{S}_{V,c}^n$ [ADV10,

SV04]. Consider the pair of primal and dual cone linear programs (LPs)

$$\begin{array}{ll}
\text{minimize} & c^T x \\
\text{subject to} & \sum_{i=1}^m x_i A_i + X = B \\
& X \in \mathbf{S}_{V,+}^n
\end{array}
\qquad
\begin{array}{ll}
\text{maximize} & -B \bullet S \\
\text{subject to} & A_i \bullet S + c_i = 0, \quad i = 1, \dots, m \\
& S \in \mathbf{S}_{V,c}^n
\end{array}
\tag{7}$$

with variables $x \in \mathbf{R}^m$, $X, S \in \mathbf{S}_V^n$, and problem parameters $A_i, B \in \mathbf{S}_V^n$, $c \in \mathbf{R}^m$. Cone LPs of this type have been studied in sparse semidefinite programming with the goal of exploiting aggregate sparsity in the coefficient matrices A_i and B [SV04, Bur03, ADV10]. Matrix completion techniques and chordal sparse matrix properties were first applied to semidefinite programming algorithms by Fukuda *et al.* [FKMN00] in a sparse implementation of the HRVW/KSH/M primal-dual algorithm.

A primal barrier method for (7) requires at each iteration the evaluation of the gradient $\nabla f(X)$ and the solution of a positive definite equation $H\Delta y = g$ with coefficients $H_{ij} = A_i \bullet (\nabla^2 f(X)[A_j])$. Efficient techniques for evaluating the Hessian $\nabla^2 f(X)[A_j]$ are therefore important in large-scale implementations. A dual barrier method for the cone programs involves evaluations of the gradient $\nabla f_*(S)$ and a set of linear equations $H\Delta y = g$ with coefficients $H_{ij} = A_i \bullet (\nabla^2 f_*(S)^{-1}[A_j])$. From the relations (5), we see that this requires the inverse \hat{X} of the maximum determinant positive definite matrix completion of S and the evaluation of the Hessian $\nabla^2 f(\hat{X})[A_j]$ at \hat{X} . We refer the reader to [ADV10, SV04] for more details.

The problem of computing a projected inverse $\mathcal{P}(X^{-1})$ (and, more generally, computing a subset of the entries of the inverse of a sparse positive definite matrix) has also been studied in statistics [GP80, ADR⁺10]. Efficient algorithms for computing the gradient of f are important in maximum likelihood estimation problems involving Gaussian distributions, for example, in sparse inverse covariance selection [DVR08]. Consider, for example, the covariance selection problem with 1-norm penalty

$$\text{minimize} \quad C \bullet X - \log \det X + \rho \|X\|_1,$$

which has been studied by several authors [HLPL06, BEd08, FHT08, dBE08, SMG10, LT10]. In this problem, X is the inverse covariance matrix of a Gaussian random variable and C is a sample covariance. The first two terms in the objective form the negative log-likelihood function of X (up to constants), and the penalty term is added to promote sparsity in the solution X . In problems of high dimension, it may be unrealistic and impractical to regard X as a dense matrix variable. Instead, one can start with a partially specified pattern based on prior knowledge, and use the penalized covariance selection to identify additional zeros. The problem can then be posed as an optimization problem over \mathbf{S}_V^n , where V is the known sparsity pattern. This greatly simplifies the cost of calculating the gradient of the smooth terms in the objective and makes it possible to solve very large covariance selection problems using first-order methods that require the gradient of $-\log \det X$ at each iteration.

Other applications include matrix approximation problems with sparse positive definite matrices as variables, for example, computing sparse quasi-Newton updates [Fle95, Yam08].

1.4 Related work and outline of the paper

We refer to the algorithms in this paper as *multifrontal* and *supernodal* because of their resemblance to multifrontal and supernodal multifrontal algorithms for Cholesky factorization [DR83, Liu92].

The multifrontal Cholesky factorization is reviewed in Section 3 and a supernodal variant, formulated in terms of clique trees, is described in Section 7.

In Section 4, we introduce multifrontal algorithms for computing the gradients of f and f_* . Similar algorithms for evaluating ∇f are discussed in [CD95, ADR⁺10]. The close connection between the problem of computing the gradient of $f(X) = -\log \det X$ and a Cholesky factorization is easily understood from the chain rule of differentiation. A practical method for computing $f(X)$ will calculate a sparse Cholesky factorization of X , for example, $X = LDL^T$ with L unit lower triangular and D diagonal, and then evaluate $f(X) = -\sum_i \log D_{ii}$. By applying the chain rule to a sparse factorization algorithm, the gradient of $f(X)$ can be evaluated at essentially the cost of the factorization itself. Moreover, the differentiation can be automated using reverse automatic differentiation software [GW08]. Although the algorithm in Section 4.1 can be obtained from the chain rule, we give a straightforward direct derivation. This not only simplifies the notation and description of the algorithm, it also helps reduce the memory requirements, which can be high in a straightforward application of reverse differentiation because of the large number of intermediate auxiliary matrices (update and frontal matrices) generated during the factorization.

As mentioned earlier, evaluating the gradient $\nabla f_*(S)$ is equivalent to inverting the mapping $X \mapsto \nabla f(X)$, and an algorithm for evaluating ∇f_* is therefore easily derived from the algorithm for ∇f (see Section 4.2).

In Section 5, we examine the problem of computing the Hessians and inverse Hessians of f and f_* , and more specifically, the problem of evaluating expressions of the form $\nabla^2 f(X)[U]$ and $\nabla^2 f_*(S)[V]$. Again, the algorithms follow conceptually from the chain rule and can be obtained by applying automatic differentiation techniques. An explicit description allows us to optimize the efficiency and memory requirements. This is particularly important in the case of sparse arguments U, V . As an important by-product, we define a factorization $\nabla^2 f(X) = \mathcal{R}^{\text{adj}} \circ \mathcal{R}$ and present efficient methods for evaluating the factors \mathcal{R} and \mathcal{R}^{adj} separately.

The methods presented in the paper are closely related to the barrier evaluation algorithms of [DVR08, DV09]. These algorithms were formulated as recursions over clique trees and can be interpreted as supernodal versions of the multifrontal algorithms presented in Sections 3–5. We elaborate on the connections in Sections 6 and 7. In contrast to the clique tree methods in [DVR08, DV09], the algorithms described here work with data structures that are widely used in sparse Cholesky factorization algorithms (namely, the compressed column storage format and elimination trees). As a result they are simpler to implement and more readily combined with techniques from the recent literature on sparse matrix factorization algorithms. Some possible further improvements in this direction are mentioned in the conclusions (Section 8).

1.5 Notation

We identify a symmetric sparsity pattern V with the positions of its lower-triangular nonzeros. In other words, a symmetric matrix X has sparsity pattern V if $X_{ij} = X_{ji} = 0$ for $(i, j) \notin V$. The entries X_{ij} and X_{ji} for $(i, j) \in V$ are treated as (structurally) nonzero, although they are allowed to be numerically zero. A lower-triangular matrix L has sparsity pattern V if the symmetric matrix $L + L^T$ has sparsity pattern V .

An *index set* is a sorted subset of the integers $\{1, 2, \dots, n\}$. The number of elements in the index set I is denoted $|I|$ and its k th element $I(k)$. If I and J are two index sets with $J \subset I$, we

define an $|I| \times |J|$ matrix E_{IJ} with entries

$$(E_{IJ})_{ij} = \begin{cases} 1 & I(i) = J(j) \\ 0 & \text{otherwise.} \end{cases}$$

This notation will be used in expressions $E_{IJ}^T B E_{IJ}$ and $E_{IJ} C E_{IJ}^T$, which have the following meaning: if A is a symmetric matrix of order n and $B = A_{II}$ is the principal submatrix indexed by I , then the matrix $E_{IJ}^T B E_{IJ}$ is equal to A_{JJ} , the principal submatrix indexed by J . The adjoint operation $B = E_{IJ} C E_{IJ}^T$, applied to a symmetric matrix C of order $|J|$, can be interpreted as first embedding C as the $|J| \times |J|$ -block A_{JJ} of an otherwise zero $n \times n$ matrix A , and then extracting the $|I| \times |I|$ submatrix $B = A_{II}$.

2 Elimination trees

This section provides some background on sparse matrices and elimination trees [Liu90, Dav06]. We define a Cholesky factorization as a factorization

$$X = LDL^T$$

with L unit lower-triangular and D positive diagonal. The sparsity pattern V of the Cholesky factor L has the following fundamental property:

$$i > j > k, \quad (i, k) \in V, \quad (j, k) \in V \implies (i, j) \in V. \quad (8)$$

(In other words, excluding accidental cancellation, $L_{ik} \neq 0$ and $L_{jk} \neq 0$ implies $L_{ij} \neq 0$.) This property distinguishes a Cholesky factor from a general sparse lower-triangular matrix. In the example in Figure 1, the presence of nonzeros in positions (5, 3) and (15, 3) implies that the entry (15, 5) is nonzero. The nonzeros in positions (9, 5), (15, 5), and (16, 5) imply that the entries in positions (15, 9), (16, 9), and (16, 15) are nonzero.

We use the notation I_k to denote the sorted set of row indices of the nonzero entries below the diagonal in column k of L . We also define $J_k = I_k \cup \{k\}$. The property (8) implies that I_k defines a complete subgraph of the filled graph, *i.e.*, the matrix $L_{J_k J_k}$ is a dense lower-triangular matrix. For example, it can be verified that the submatrix indexed by $J_5 = \{5, 9, 15, 16\}$ in Figure 1 is dense. The number of nonzeros below the diagonal in column k , *i.e.*, the cardinality $|I_k|$ of I_k , is called the *monotone degree* of vertex k .

The *elimination tree* (etree) is defined in terms of the sparsity pattern of the factor L as follows. It is a tree (or a forest if L is reducible) with n vertices, labeled 1 to n . The parent of vertex k is the row index j of the first nonzero below the diagonal of column k of L , *i.e.*, the vertex $\min I_k$. As a consequence, each vertex has a lower index than its parent, so the vertices in the elimination tree are numbered in a *topological ordering*. An example is shown in Figure 1.

We will use two important properties of elimination trees.

Theorem 1 [Liu90, theorem 3.1] *If $j \in I_k$, then j is an ancestor of k in the elimination tree.*

Note that the converse does not hold.

Theorem 2 [Liu92, theorem 3.1] *If vertex j is an ancestor of vertex k in the elimination tree, then the nonzero structure of $(L_{jk}, L_{j+1,k}, \dots, L_{nk})$ is contained in the structure of $(L_{jj}, L_{j+1,j}, \dots, L_{nj})$.*

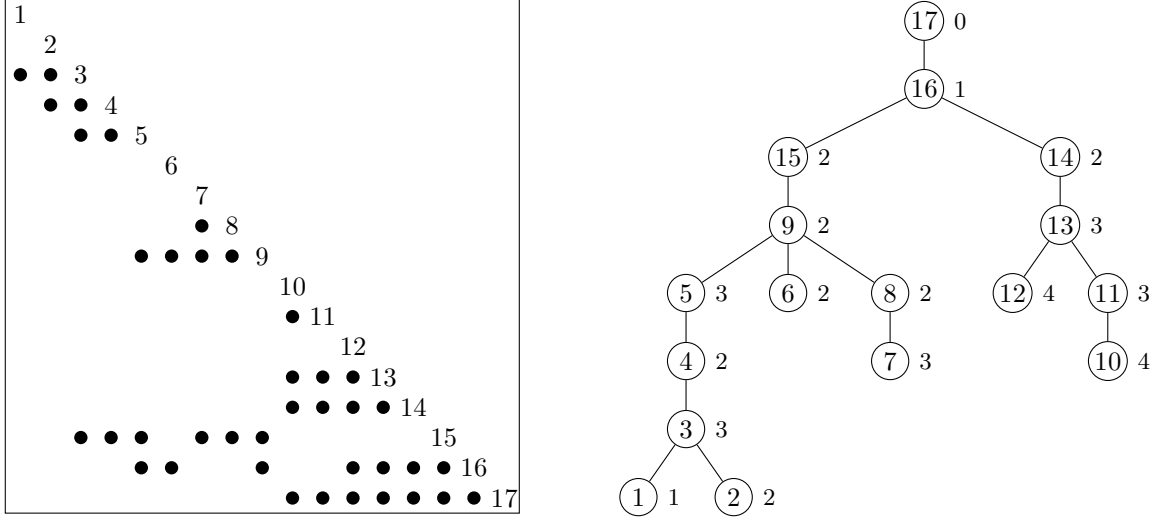


Figure 1: Filled pattern and elimination tree. The numbers next to vertices in the elimination tree are the monotone degrees $|I_k|$.

In the notation for the column structure introduced above, this theorem asserts that if j is an ancestor of k , then

$$I_k \cap \{j, j+1, \dots, n\} \subseteq J_j.$$

In particular, if j is the parent of k (hence, by definition, j is the first element of I_k and therefore $I_k \subseteq \{j, j+1, \dots, n\}$), then

$$I_k \subseteq J_j, \quad |I_k| \leq |I_j| + 1. \quad (9)$$

By applying these inequalities recursively to a path i_1, i_2, \dots, i_r from a vertex i_1 in the elimination tree to one of its ancestors i_r , we obtain a chain of inclusions

$$I_{i_1} \subseteq J_{i_2} \subseteq \{i_2\} \cup J_{i_3} \subseteq \dots \subseteq \{i_2, i_3, \dots, i_{r-1}\} \cup J_{i_r} \quad (10)$$

and inequalities

$$|I_{i_1}| \leq |I_{i_2}| + 1 \leq |I_{i_3}| + 2 \leq \dots \leq |I_{i_r}| + r - 1. \quad (11)$$

This can be verified in the elimination tree in Figure 1. As we move along a path from a leaf vertex to the root of the tree, the monotone degrees can increase or decrease, but they never decrease by more than one per step.

3 Cholesky factorization and multiplication

In this section, we review the multifrontal algorithm for Cholesky factorization [DR83, Liu92]. We then describe a similar algorithm for the related problem of computing a matrix, given its Cholesky factors.

3.1 Cholesky factorization

Recall that we define the Cholesky factorization as a decomposition $X = LDL^T$, with D positive diagonal and L unit lower-triangular. The formulas for L and D are easily derived from the equation

$X = LDL^T$. The $J_j \times J_j$ block of the factorization is

$$\begin{bmatrix} X_{jj} & X_{I_j j}^T \\ X_{I_j j} & X_{I_j I_j} \end{bmatrix} = \sum_{k < j} D_{kk} \begin{bmatrix} L_{jk} \\ L_{I_j k} \end{bmatrix} \begin{bmatrix} L_{jk} \\ L_{I_j k} \end{bmatrix}^T + D_{jj} \begin{bmatrix} 1 \\ L_{I_j j} \end{bmatrix} \begin{bmatrix} 1 \\ L_{I_j j} \end{bmatrix}^T + \sum_{k > j} D_{kk} \begin{bmatrix} 0 \\ L_{I_j k} \end{bmatrix} \begin{bmatrix} 0 \\ L_{I_j k} \end{bmatrix}^T.$$

The first column of the equation is

$$\begin{bmatrix} X_{jj} \\ X_{I_j j} \end{bmatrix} = \sum_{k < j} D_{kk} L_{jk} \begin{bmatrix} L_{jk} \\ L_{I_j k} \end{bmatrix} + D_{jj} \begin{bmatrix} 1 \\ L_{I_j j} \end{bmatrix}. \quad (12)$$

The multifrontal algorithm takes advantage of properties of the elimination tree associated with L to compute the sum on the right-hand side. First, we recall that $L_{jk} \neq 0$ only if k is a descendant of j in the elimination tree (Theorem 1). The sum in (12) can therefore be replaced by a sum over the proper descendants of j . The set of proper descendants of vertex j is

$$T_j \setminus \{j\} = \bigcup_{i \in \text{ch}(j)} T_i,$$

where T_j is the subtree of the elimination tree rooted at vertex j and $\text{ch}(j)$ are the children of vertex j . The equation (12) then becomes

$$\begin{bmatrix} X_{jj} \\ X_{I_j j} \end{bmatrix} = \sum_{i \in \text{ch}(j)} \sum_{k \in T_i} D_{kk} L_{jk} \begin{bmatrix} L_{jk} \\ L_{I_j k} \end{bmatrix} + D_{jj} \begin{bmatrix} 1 \\ L_{I_j j} \end{bmatrix}. \quad (13)$$

Second, suppose that for each vertex i in the elimination tree, we define a dense matrix

$$U_i = - \sum_{k \in T_i} D_{kk} L_{I_i k} L_{I_i k}^T. \quad (14)$$

The matrix U_i is called the *update matrix* for vertex i . Using the definition of E_{IJ} in Section 1.5, we can write

$$- \sum_{k \in T_i} D_{kk} \begin{bmatrix} L_{jk} \\ L_{I_j k} \end{bmatrix} \begin{bmatrix} L_{jk} \\ L_{I_j k} \end{bmatrix}^T = E_{J_j I_i} U_i E_{J_j I_i}^T. \quad (15)$$

This follows from Theorem 2: if $i \in \text{ch}(j)$ and $k \in T_i$, then

$$I_k \cap \{j, j+1, \dots, n\} \subseteq I_i \subseteq J_j$$

and therefore $L_{J_j k} = E_{J_j I_i} L_{I_i k}$. (The multiplication with $E_{J_j I_i}$ copies the entries $L_{I_i k}$ to the correct position in $L_{J_j k}$ and inserts zeros for the other entries.) Adding the first columns of $E_{J_j I_i} U_i E_{J_j I_i}^T$ for all $i \in \text{ch}(j)$ therefore gives the first term on the right-hand side of (13).

Thus, by combining the nonzero lower-triangular entries in column j of X (the left-hand side of (13)) and the update matrices of the children of vertex j (to assemble the sum on the right-hand side), we collect all the information needed to compute D_{jj} and $L_{I_j j}$ from (13).

Furthermore, the same equation (13) shows how the update matrix U_j for vertex j can be calculated. This is clearer if we rewrite (13) in matrix form using (15) as

$$\begin{bmatrix} X_{jj} & X_{I_{jj}}^T \\ X_{I_{jj}} & 0 \end{bmatrix} + \sum_{i \in \text{ch}(j)} E_{J_j I_i} U_i E_{J_j I_i}^T = D_{jj} \begin{bmatrix} 1 \\ L_{I_{jj}} \end{bmatrix} \begin{bmatrix} 1 \\ L_{I_{jj}} \end{bmatrix}^T + \begin{bmatrix} 0 & 0 \\ 0 & U_j \end{bmatrix} \quad (16)$$

$$= \begin{bmatrix} 1 & 0 \\ L_{I_{jj}} & I \end{bmatrix} \begin{bmatrix} D_{jj} & 0 \\ 0 & U_j \end{bmatrix} \begin{bmatrix} 1 & L_{I_{jj}}^T \\ 0 & I \end{bmatrix}. \quad (17)$$

The first column of this equation is identical to (13). The 2,2 block follows from the definition of U_j and the identity (15). The matrix on the left-hand side of (17) is called the j th *frontal matrix*. The equation (17) shows that once the frontal matrix has been assembled, we can compute D_{jj} , $L_{I_{jj}}$, and U_j by a pivot step.

The resulting algorithm to compute L , D , given a positive definite X , is summarized below.

Algorithm 3.1. *Cholesky factorization.*

Input. A positive definite matrix X .

Output. The factors L , D in the Cholesky factorization $X = LDL^T$.

Algorithm. Iterate over $j \in \{1, \dots, n\}$ in topological order (*i.e.*, visiting each vertex of the elimination tree before its parent). For each j , form the frontal matrix

$$F_j = \begin{bmatrix} F_{11} & F_{21}^T \\ F_{21} & F_{22} \end{bmatrix} = \begin{bmatrix} X_{jj} & X_{I_{jj}}^T \\ X_{I_{jj}} & 0 \end{bmatrix} + \sum_{i \in \text{ch}(j)} E_{J_j I_i} U_i E_{J_j I_i}^T \quad (18)$$

and calculate D_{jj} , the j th column of L , and the j th update matrix U_j from

$$D_{jj} = F_{11}, \quad L_{I_{jj}} = \frac{1}{D_{jj}} F_{21}, \quad U_j = F_{22} - D_{jj} L_{I_{jj}} L_{I_{jj}}^T. \quad (19)$$

In a practical implementation, with the lower-triangular part of X stored in a sparse format (typically, the compressed column structure or CCS; see [Dav06]), one can overwrite X_{jj} with D_{jj} and $X_{I_{jj}}$ with $L_{I_{jj}}$ after cycle j . The auxiliary matrices F_j and U_i are stored as dense matrices (either as two separate arrays or by letting U_j overwrite the 2, 2 block of F_j). The main step in the algorithm is the level-2 BLAS operation in the calculation of U_j in (19) [DCHH88]. The frontal matrix F_j can be discarded after the vertex j has been processed. The update matrix U_j can be discarded after the parent of vertex j has been processed.

3.2 Cholesky multiplication

The equation (16) also shows how the j th column of X can be computed from D_{jj} , column j of L , and the update matrices for the children of vertex j . This yields an algorithm for the inverse operation of the Cholesky factorization, *i.e.*, the matrix multiplication LDL^T , which will be important in Section 4.

Algorithm 3.2. *Cholesky product.*

Input. Cholesky factors L, D .

Output. The matrix $X = LDL^T$.

Algorithm. Iterate over $j \in \{1, \dots, n\}$ in topological order. For each j , calculate X_{jj} , $X_{I_{jj}}$, and U_j from

$$\begin{bmatrix} X_{jj} & X_{I_{jj}}^T \\ X_{I_{jj}} & -U_j \end{bmatrix} = D_{jj} \begin{bmatrix} 1 \\ L_{I_{jj}} \end{bmatrix} \begin{bmatrix} 1 \\ L_{I_{jj}} \end{bmatrix}^T - \sum_{i \in \text{ch}(j)} E_{J_j I_i} U_i E_{J_j I_i}^T. \quad (20)$$

The matrices U_i can be deleted after the parent of vertex i has been processed. In a practical implementation, we compute the left-hand side of (20) as a dense matrix, via a level-2 BLAS operation for the outer-product on the right-hand side. Then X_{jj} and $X_{I_{jj}}$ are copied to the CCS structure for X .

4 Gradients

In this section we describe ‘multifrontal’ algorithms for evaluating the gradients of the barrier and the dual barrier.

Recall that the primal gradient is defined as $\nabla f(x) = -\mathcal{P}(X^{-1})$, where \mathcal{P} denotes projection on the filled pattern V of X . It is straightforward to show that the entries of the projected inverse $\mathcal{P}(X^{-1})$ can be computed directly from the Cholesky factors L, D without calculating any entries of X^{-1} outside V (see Section 4.1). This observation is the basis of several algorithms published in the literature. The algorithm we describe here is equivalent to the *inverse multifrontal algorithm* in [CD95], but we give a different and shorter derivation. The projected inverse algorithm in [DVR08] can be viewed as a supernodal variant of the algorithm discussed here (see the discussion in Section 7). Another closely related algorithm is described by Amestoy *et al.* [ADR⁺10], who consider the problem of computing a few entries of the inverse of a large sparse matrix.

Evaluation of the dual gradient corresponds to the inverse operation, *i.e.*, the problem of solving the nonlinear equation $\mathcal{P}(X^{-1}) = S$ with variable $X \in \mathbf{S}_V^n$. A multifrontal algorithm for this problem is derived in Section 4.2.

4.1 Primal gradient

Define $Z = X^{-1}$ and $S = \mathcal{P}(Z)$. We are interested in an efficient method for computing S from the Cholesky factors L and D of X . The matrix $Z = L^{-T} D^{-1} L^{-1}$ satisfies

$$ZL = L^{-T} D^{-1}.$$

The $J_j \times j$ block of this equation only involves entries of Z in the projection $S = \mathcal{P}(Z)$:

$$\begin{bmatrix} S_{jj} & S_{I_{jj}}^T \\ S_{I_{jj}} & S_{I_{jj} I_j} \end{bmatrix} \begin{bmatrix} 1 \\ L_{I_{jj}} \end{bmatrix} = \begin{bmatrix} 1/D_{jj} \\ 0 \end{bmatrix}. \quad (21)$$

The vertices of I_j are ancestors of vertex j (Theorem 1). Therefore, if we calculate the columns of S following a reverse topological order of the vertices of the elimination tree, then the matrix $S_{I_j I_j}$ is known when we arrive at column j . Given $S_{I_j I_j}$, it is easy to compute $S_{I_j j}$ and S_{jj} from (21):

$$S_{I_j j} = -S_{I_j I_j} L_{I_j j}, \quad S_{jj} = \frac{1}{D_{jj}} - S_{I_j j}^T L_{I_j j}. \quad (22)$$

Accessing the vectors $L_{I_j j}$ and $S_{I_j j}$ is easy if L and the lower-triangular part of S are stored in a CCS data structure. Retrieving $S_{I_j I_j}$ from the CCS representation of S can be avoided by using an idea similar to the multifrontal Cholesky algorithm. For each vertex j of the elimination tree, we define a dense ‘update matrix’

$$V_j = S_{I_j I_j}.$$

It follows from the properties of the elimination tree (theorem 2 and equation (9)) and the definition of $E_{J_j I_i}$ that if i is a child of vertex j , then

$$V_i = E_{J_j I_i}^T \begin{bmatrix} S_{jj} & S_{I_j j}^T \\ S_{I_j j} & V_j \end{bmatrix} E_{J_j I_i}.$$

By using this formula to propagate V_i , we obtain a ‘multifrontal’ algorithm for computing $\mathcal{P}(X^{-1})$.

Algorithm 4.1. *Projected inverse.*

Input. The Cholesky factors L , D of a positive definite matrix X .

Output. The projected inverse $S = \mathcal{P}(X^{-1}) = -\nabla f(X)$.

Algorithm. Iterate over $j = \{1, 2, \dots, n\}$ in reverse topological order (*i.e.*, visiting each vertex before its children). For each j , calculate S_{jj} and $S_{I_j j}$ from

$$S_{I_j j} = -V_j L_{I_j j}, \quad S_{jj} = \frac{1}{D_{jj}} - S_{I_j j}^T L_{I_j j}, \quad (23)$$

and compute the update matrices

$$V_i = E_{J_j I_i}^T \begin{bmatrix} S_{jj} & S_{I_j j}^T \\ S_{I_j j} & V_j \end{bmatrix} E_{J_j I_i}, \quad i \in \text{ch}(j). \quad (24)$$

The matrix V_j can be discarded after cycle j , and $S_{I_j j}$ and S_{jj} can overwrite $L_{I_j j}$ and D_{jj} in a CCS data structure. As in Algorithm 3.1, the algorithm involves operations with the dense matrices V_j and

$$\begin{bmatrix} S_{jj} & S_{I_j j}^T \\ S_{I_j j} & S_{I_j I_j} \end{bmatrix}.$$

The main calculation is a level-2 BLAS operation (the matrix-vector product $V_j L_{I_j j}$).

4.2 Dual gradient

As mentioned in the introduction, the solution $X \in \mathbf{S}_V^n$ of the equation $\mathcal{P}(X^{-1}) = S$ is the inverse of the maximum determinant positive definite completion of S , and it is also the negative of $\nabla f_*(S)$. The Cholesky factors of X can be computed by solving for $L_{I_j j}$ and D_{jj} from (22) as in the following algorithm.

Algorithm 4.2. *Matrix completion.*

Input. A matrix $S \in \mathbf{S}_V^n$ that has a positive definite completion.

Output. The Cholesky factors L, D of $X = -\nabla f_*(S)$, i.e., the positive definite matrix $X \in \mathbf{S}_V^n$ that satisfies $\mathcal{P}(X^{-1}) = S$.

Algorithm. Iterate over $j \in \{1, 2, \dots, n\}$ in reverse topological order. For each j , compute D_{jj} and the j th column of L from

$$L_{I_j j} = -V_j^{-1} S_{I_j j}, \quad D_{jj} = (S_{jj} + S_{I_j j}^T L_{I_j j})^{-1}, \quad (25)$$

and compute the update matrices

$$V_i = E_{J_j I_i}^T \begin{bmatrix} S_{jj} & S_{I_j j}^T \\ S_{I_j j} & V_j \end{bmatrix} E_{J_j I_i}, \quad i \in \text{ch}(j). \quad (26)$$

The update matrix V_j can be discarded after cycle j , and $L_{I_j j}$ and D_{jj} can overwrite $S_{I_j j}$ and S_{jj} in a CCS data structure.

The cost of Algorithm 4.2 is higher than that of Algorithm 4.1 because step (25) involves the solution of an equation with V_j as coefficient matrix, whereas (24) only requires a multiplication. It is therefore of interest to propagate a factorization of V_j instead of the matrix itself, and to replace (26) by an efficient method for computing a factorization of V_i , given a factorization of V_j . This idea can be implemented as follows. We use a factorization of the form $V_j = R_j R_j^T$, with R_j upper triangular of order $|I_j|$. We need to replace (26) with an efficient method for computing R_i from R_j . The matrix $E_{J_j I_i}$ can be partitioned as

$$E_{J_j I_i} = \begin{bmatrix} 1 & 0 \\ 0 & E \end{bmatrix},$$

where $E = E_{I_j, I_i \setminus \{j\}}$. (This follows from the fact that j is the first element of J_j by definition of J_j , and j is also the first element of I_i if j is the parent of vertex i . As a consequence, the 1,1 element of matrix $E_{J_j I_i}$ is equal to one). From (26),

$$V_i = E_{J_j I_i}^T \begin{bmatrix} S_{jj} & S_{I_j j}^T \\ S_{I_j j} & V_j \end{bmatrix} E_{J_j I_i} = \begin{bmatrix} a & b^T \\ b & C C^T \end{bmatrix}$$

where

$$a = S_{jj}, \quad b = E^T S_{I_j j}, \quad C = E^T R_j.$$

The matrix C is obtained from the upper triangular matrix R_j by deleting the rows in $J_j \setminus I_i$. It can be reduced to square upper triangular form by writing it as $C = R Q^T$ with R upper triangular

and Q orthogonal (a product of Householder transformations [GV96]). Then the triangular factor in $V_i = R_i R_i^T$ is given by

$$R_i = \begin{bmatrix} (a - \|R^{-1}b\|_2^2)^{1/2} & (R^{-1}b)^T \\ 0 & R \end{bmatrix}.$$

This is summarized below.

Algorithm 4.3. *Matrix completion with factored update matrices.*

Input. A matrix $S \in \mathbf{S}_V^n$ that has a positive definite completion.

Output. The Cholesky factors L, D of $X = -\nabla f_*(S)$, i.e., of the positive definite matrix $X \in \mathbf{S}_V^n$ that satisfies $\mathcal{P}(X^{-1}) = S$.

Algorithm. Iterate over $j \in \{1, 2, \dots, n\}$ in reverse topological order. For each j ,

- compute D_{jj} and the j th column of L from

$$L_{I_{jj}} = -R_j^{-T} R_j^{-1} S_{I_{jj}}, \quad D_{jj} = (S_{jj} + S_{I_{jj}}^T L_{I_{jj}})^{-1} \quad (27)$$

- for $i \in \text{ch}(j)$, compute a factorization

$$E_{I_j, I_i \setminus \{j\}}^T R_j = R Q^T \quad (28)$$

with R upper triangular and Q orthogonal, and compute

$$R_i = \begin{bmatrix} (S_{jj} - \|R^{-1} S_{I_{jj}}\|_2^2)^{1/2} & (R^{-1} S_{I_{jj}})^T \\ 0 & R \end{bmatrix}. \quad (29)$$

The cost of step (27) is order $|I_j|^2$ for the forward and back substitutions, and the cost of (29) is proportional to $|I_i|^2$ (for the computation of $R^{-1} S_{I_{jj}}$). The cost of the reduction to triangular form in (28) is difficult to quantify because it depends on the number of rows in R_j that are deleted in the multiplication $E^T R_j$ and on their positions. However the total cost is usually much less than the cost of computing R_i from scratch, as shown by the experiments in the next section.

4.3 Numerical results

In this section, we give experimental results with Algorithms 4.1, 4.2, and 4.3. The algorithms were implemented in Python 2.7, using the Python library CVXOPT version 1.1.3 [DV10] and its interfaces to LAPACK and BLAS¹ for the sparse and dense matrix computations. Some critical code segments were implemented in C (such as the Householder updates in Algorithm 4.3 and the “extend-add” operation $F := F + E_{JI} U E_{JI}^T$ and its adjoint.) The experiments were conducted on an Intel Q6600 CPU (2.4 GHz Core 2 Quad) computer with 4 GB memory, running Ubuntu 11.04.

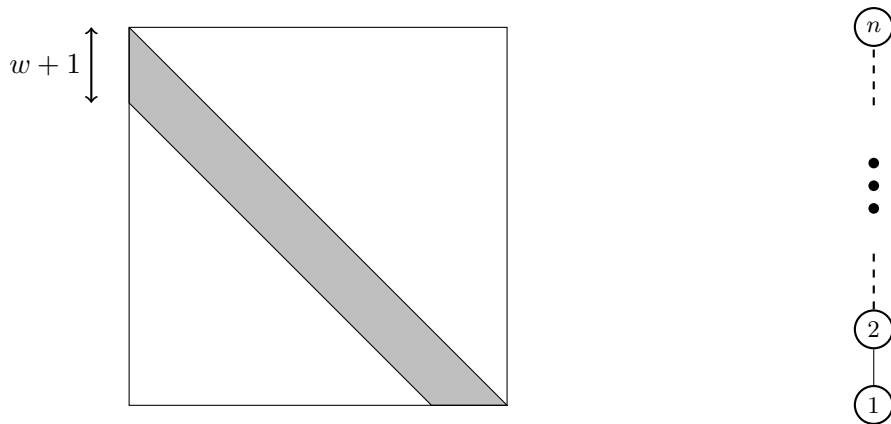


Figure 2: Lower-triangular band pattern with bandwidth w and the corresponding elimination tree. The vertices $1, \dots, n-w$ have monotone degree w . The vertices $k = n-w+1, \dots, n$ have monotone degree $n-k$.

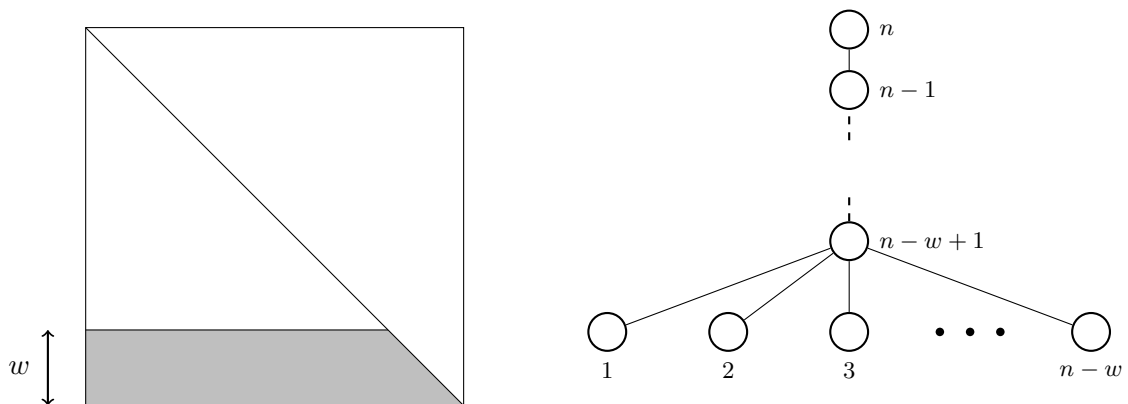


Figure 3: Lower-triangular arrow pattern with width w and the corresponding elimination tree. The vertices $1, \dots, n-w$ have monotone degree w . The vertices $k = n-w+1, \dots, n$ have monotone degree $n-k$.

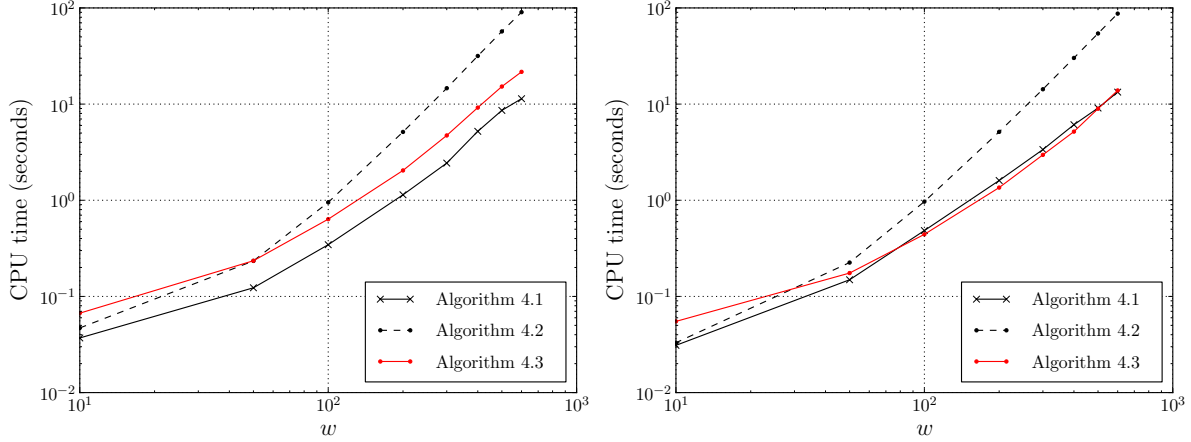


Figure 4: Cost of computing the primal gradient (Algorithm 4.1) and the dual gradient (Algorithms 4.2 and 4.3) for (a) band patterns and (b) arrow patterns of order $n = 2000$, as a function of the pattern width w .

4.3.1 Band and arrow patterns

Band and arrow patterns are two basic sparsity patterns for which the complexity of the algorithms is easy to analyze. We assume $n \gg w$, where w is the bandwidth or blockwidth (see Figures 2 and 3). It is easy to see that the cost of a Cholesky factorization of a matrix with one of these two patterns is $O(nw^2)$.

Step (23) of Algorithm 4.1 involves matrix-vector multiplications of order $|I_j|$. The complexity of the algorithm is dominated by the total cost of these products, *i.e.*, $O(nw^2)$. This is similar to the cost of a Cholesky factorization. Step (25) of Algorithm 4.2 on the other hand requires solving a dense positive definite system of order $|I_j|$. The total complexity is therefore $O(nw^3)$.

In Algorithm 4.3, the cost of step (27) is reduced to $O(w^2)$ per iteration. For $j = n - w + 1, \dots, n$, we have $E = I$ in step (24), so C is upper triangular and we only need to compute $R^{-1}b$. For the other vertices in the elimination tree ($j = 1, \dots, n - w$), $C = E^T R_j$ is the matrix R_j with one row deleted: the last row in the case of a band pattern, the first row in the case of an arrow pattern. The cost of reducing C to triangular form is therefore zero in the case of an arrow pattern and $O(w^2)$ in the case of a band pattern. In either case, the total cost of Algorithm 4.3 is reduced to $O(nw^2)$.

In Figure 4, we compare the CPU times of the three algorithms as a function of the width w . As can be seen, the cost of computing the dual gradient using Algorithm 4.3 is comparable to the cost of the primal gradient using Algorithm 4.1, and the cost of the two algorithms grows roughly as w^2 for fixed n . Notice the small gap between the cost of Algorithms 4.3 and 4.1 for band patterns. This gap reflects the cost of reducing C to triangular form in Algorithm 4.3.

4.3.2 General sparse patterns

In the second experiment, we use a benchmark set of large symmetric sparsity patterns from the University of Florida Sparse Matrix Collection [Dav09]. The AMD ordering was used to compute

¹We link against the single-threaded reference implementations of BLAS and LAPACK in Ubuntu.

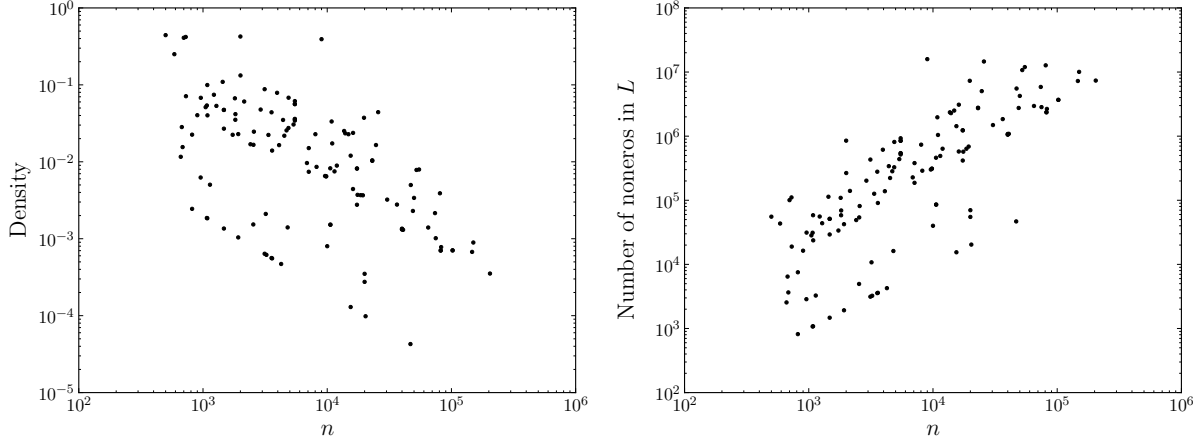


Figure 5: Scatter plot of n versus (a) the density ($|V|/(n(n+1)/2)$) and (b) the number of nonzeros in L for the 128 test problems. Each dot represents a problem from the University of Florida Sparse Matrix collection.

filled patterns. To prevent out-of-core computations, we restrict the experiment to matrices for which the filled pattern occupies less than 250 MB of memory. The set of test problems includes 128 sparsity patterns, with n ranging from 500 to 204316 and with $|V|$ between 817 and 15,894,180. A scatter plot of the number of nonzeros and the density of the test problems versus the dimension n is shown in Figure 5.

Figure 6 shows the CPU times for Algorithm 4.1 (primal gradient or projected inverse) and 4.3 (dual gradient or completion), and for Algorithms 4.3 and 4.2 (completion with and without Householder updates, respectively). Each dot represents one of the sparsity patterns in the test set. The results indicate that in practice, on this set of realistic sparsity patterns, the costs of computing the primal and dual gradients are comparable.

5 Hessian

The Hessian $\mathcal{H} = \nabla^2 f(X)$ of f at X is defined as

$$\mathcal{H}(Y) = \mathcal{P}(X^{-1}YX^{-1}) = - \left. \frac{d}{dt} \mathcal{P}(X + tY)^{-1} \right|_{t=0}.$$

A method for evaluating $\mathcal{H}(Y)$ can therefore be found by differentiating the formulas for evaluating the gradient $-\mathcal{P}(X^{-1})$. As we have seen, $\mathcal{P}(X^{-1})$ is obtained in two stages. First the Cholesky factors L, D of X are computed, column by column, following a topological ordering of the elimination tree (Algorithm 3.1). Then the projected inverse $\mathcal{P}(X^{-1})$ is computed from L and D , column by column, in reverse topological order (Algorithm 4.1). Linearizing the two algorithms will provide an algorithm for $\mathcal{H}(Y)$. We give the details in Section 5.2.

We also consider the problem of evaluating $Y = \mathcal{H}^{-1}(T)$, *i.e.*, solving the linear equation

$$\mathcal{P}(X^{-1}YX^{-1}) = T$$

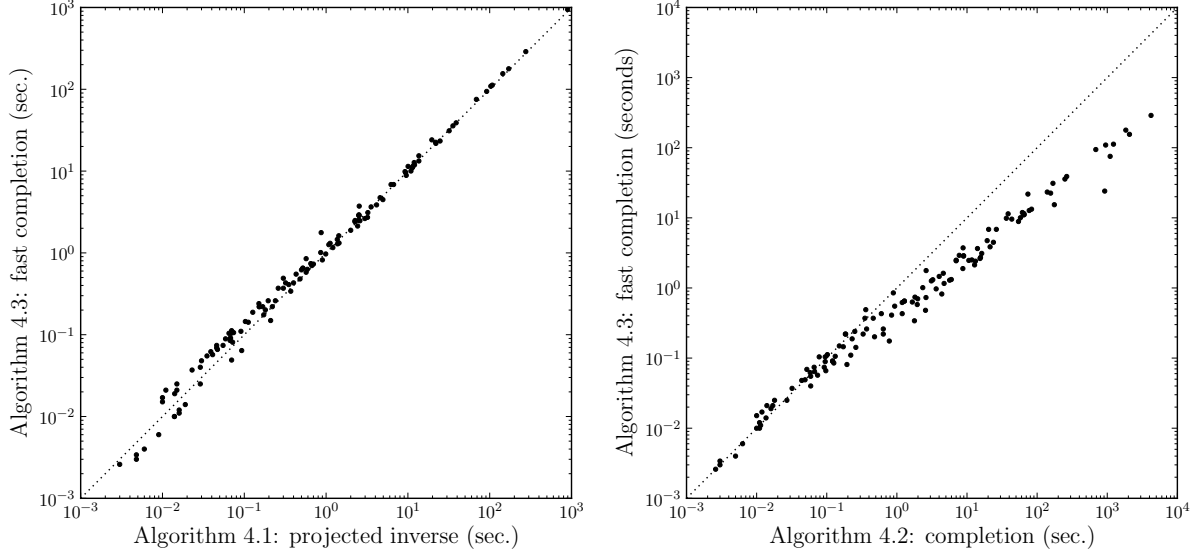


Figure 6: CPU times for primal and dual gradient algorithms for a test set of 128 sparsity patterns: (a) shows the cost of Algorithm 4.1 (primal gradient or projected inverse) versus the cost of Algorithm 4.3 (dual gradient or matrix completion) for different sparsity patterns; (b) compares the cost of computing the dual gradient via Algorithms 4.3 and 4.2 *i.e.*, with and without the update technique described in Section 4.2.

for $Y \in \mathbf{S}_V^n$, given $T \in \mathbf{S}_V^n$. An algorithm for this problem can be formulated by inverting the calculation of $\mathcal{H}(Y)$ or, alternatively, by linearizing the algorithms for matrix completion (Algorithms 4.2 and 4.3) and the Cholesky product (Algorithm 3.2); see Section 5.3.

In applications, it is often useful to know a factorization of the Hessian as a composition of a mapping and its adjoint,

$$\mathcal{H}(Y) = \mathcal{R}^{\text{adj}}(\mathcal{R}(Y)).$$

Such a factorization is discussed in Section 5.4.

5.1 Linearized Cholesky factorization and matrix completion

Let $L(t)$, $D(t)$ be the matrices in the factorization $X(t) = X + tY = L(t)D(t)L(t)^T$ and let $U_i(t)$ be the i th update matrix in the multifrontal factorization algorithm for $X(t)$, *i.e.*,

$$U_i(t) = - \sum_{k \in T_i} D_{kk}(t) L_{I_k k}(t) L_{I_k k}(t)^T.$$

We denote by L' , D' , U'_i the derivatives of $L(t)$, $D(t)$, $U_i(t)$ at $t = 0$. These derivatives can be found by linearizing the equation (16) with X replaced by $X + tY$,

$$\begin{aligned} & \begin{bmatrix} X_{jj} & X_{I_{jj}}^T \\ X_{I_{jj}} & 0 \end{bmatrix} + t \begin{bmatrix} Y_{jj} & Y_{I_{jj}}^T \\ Y_{I_{jj}} & 0 \end{bmatrix} + \sum_{i \in \text{ch}(j)} E_{J_j I_i} U_i(t) E_{J_j I_i}^T \\ &= D_{jj}(t) \begin{bmatrix} 1 \\ L_{I_{jj}}(t) \end{bmatrix} \begin{bmatrix} 1 \\ L_{I_{jj}}(t) \end{bmatrix}^T + \begin{bmatrix} 0 & 0 \\ 0 & U_j(t) \end{bmatrix}. \end{aligned}$$

Taking the derivatives of the left- and right-hand sides at $t = 0$ gives

$$\begin{bmatrix} Y_{jj} & Y_{I_{jj}}^T \\ Y_{I_{jj}} & 0 \end{bmatrix} + \sum_{i \in \text{ch}(j)} E_{J_j I_i} U_i' E_{J_j I_i}^T = \begin{bmatrix} 1 & 0 \\ L_{I_{jj}} & I \end{bmatrix} \begin{bmatrix} D'_{jj} & (D_{jj} L'_{I_{jj}})^T \\ D_{jj} L'_{I_{jj}} & U_j' \end{bmatrix} \begin{bmatrix} 1 & L_{I_{jj}}^T \\ 0 & I \end{bmatrix}. \quad (30)$$

This will be the key equation for computing the linearized Cholesky factors D' , L' from Y , and conversely, the linearized Cholesky product Y from the linearized factors D' , L' .

Similarly, we define $Z(t) = (X + tY)^{-1}$, $S(t) = \mathcal{P}(Z(t))$, and

$$V_i(t) = S_{I_i I_i}(t).$$

We write the derivatives of $S(t)$ and $V_i(t)$ at $t = 0$ as $-T$ and $-V_i'$. Substituting $S(t)$, $L(t)$, $D(t)$ for S , L , D in (21) gives

$$\begin{bmatrix} S_{jj}(t) & S_{I_{jj}}(t)^T \\ S_{I_{jj}}(t) & S_{I_j I_j}(t) \end{bmatrix} \begin{bmatrix} 1 \\ L_{I_{jj}}(t) \end{bmatrix} = \begin{bmatrix} 1/D_{jj}(t) \\ 0 \end{bmatrix}$$

and differentiating with respect to t gives

$$-\begin{bmatrix} T_{jj} & T_{I_{jj}}^T \\ T_{I_{jj}} & V_j' \end{bmatrix} \begin{bmatrix} 1 \\ L_{I_{jj}} \end{bmatrix} + \begin{bmatrix} S_{jj} & S_{I_{jj}}^T \\ S_{I_{jj}} & S_{I_j I_j} \end{bmatrix} \begin{bmatrix} 0 \\ L'_{I_{jj}} \end{bmatrix} = \begin{bmatrix} -D'_{jj}/D_{jj}^2 \\ 0 \end{bmatrix}.$$

Using $S_{I_{jj}} = -S_{I_j I_j} L_{I_{jj}}$ (from (22)) this can be written in a more symmetric form as

$$\begin{bmatrix} T_{jj} & T_{I_{jj}}^T \\ T_{I_{jj}} & V_j' \end{bmatrix} \begin{bmatrix} 1 & 0 \\ L_{I_{jj}} & I \end{bmatrix} = \begin{bmatrix} 1 & -L_{I_{jj}}^T \\ 0 & I \end{bmatrix} \begin{bmatrix} D'_{jj}/D_{jj}^2 & (S_{I_j I_j} L'_{I_{jj}})^T \\ S_{I_j I_j} L'_{I_{jj}} & V_j' \end{bmatrix},$$

i.e.,

$$\begin{bmatrix} 1 & L_{I_{jj}}^T \\ 0 & I \end{bmatrix} \begin{bmatrix} T_{jj} & T_{I_{jj}}^T \\ T_{I_{jj}} & V_j' \end{bmatrix} \begin{bmatrix} 1 & 0 \\ L_{I_{jj}} & I \end{bmatrix} = \begin{bmatrix} D'_{jj}/D_{jj}^2 & (S_{I_j I_j} L'_{I_{jj}})^T \\ S_{I_j I_j} L'_{I_{jj}} & V_j' \end{bmatrix}. \quad (31)$$

This equation allows us to compute T given the linearized factors D' , L' , and conversely, compute D' , L' given T .

5.2 Hessian

The algorithm for computing $T = \mathcal{H}(Y)$ first computes L' , D' by the linearized Cholesky factorization, *i.e.*, from (30), and then T from L' , D' by the linearized projected inverse algorithm, *i.e.*, from (31). To simplify the notation, we define two matrices K , $M \in \mathbf{S}_V^n$ as

$$K_{jj} = D'_{jj}, \quad K_{I_{jj}} = L'_{I_{jj}} D_{jj}, \quad M_{jj} = \frac{D'_{jj}}{D_{jj}^2}, \quad M_{I_{jj}} = S_{I_j I_j} L'_{I_{jj}}, \quad (32)$$

for $j = 1, \dots, n$.

Algorithm 5.1. *Hessian evaluation.*

Input. A matrix $Y \in \mathbf{S}_V^n$, the Cholesky factors L, D of a positive definite matrix $X \in \mathbf{S}_V^n$, and the projected inverse $S = \mathcal{P}(X^{-1})$.

Output. The matrix $T = \mathcal{H}(Y) = \mathcal{P}(X^{-1}YX^{-1})$ where \mathcal{H} is the Hessian of f at X .

Algorithm.

1. Iterate over $j \in \{1, 2, \dots, n\}$ in topological order. For each j , calculate D'_{jj} , the j th column of K , and the update matrix U'_j via

$$\begin{bmatrix} K_{jj} & K_{I_{jj}}^T \\ K_{I_{jj}} & U'_j \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ -L_{I_{jj}} & I \end{bmatrix} \left(\begin{bmatrix} Y_{jj} & Y_{I_{jj}}^T \\ Y_{I_{jj}} & 0 \end{bmatrix} + \sum_{i \in \text{ch}(j)} E_{J_j I_i} U'_i E_{J_j I_i}^T \right) \begin{bmatrix} 1 & -L_{I_{jj}}^T \\ 0 & I \end{bmatrix}.$$

2. For $j \in \{1, 2, \dots, n\}$, compute column j of M via

$$M_{jj} = \frac{K_{jj}}{D_{jj}^2}, \quad M_{I_{jj}} = \frac{1}{D_{jj}} S_{I_{jj}} K_{I_{jj}}.$$

3. Iterate over $j \in \{1, 2, \dots, n\}$ in reverse topological order. For each j , calculate T_{jj} and $T_{I_{jj}}$ from

$$\begin{bmatrix} T_{jj} & T_{I_{jj}}^T \\ T_{I_{jj}} & V'_j \end{bmatrix} = \begin{bmatrix} 1 & -L_{I_{jj}}^T \\ 0 & I \end{bmatrix} \begin{bmatrix} M_{jj} & M_{I_{jj}}^T \\ M_{I_{jj}} & V'_j \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -L_{I_{jj}} & I \end{bmatrix}$$

and the update matrices V'_i for the children of vertex j via

$$V'_i = E_{J_j I_i}^T \begin{bmatrix} T_{jj} & T_{I_{jj}}^T \\ T_{I_{jj}} & V'_j \end{bmatrix} E_{J_j I_i}, \quad i \in \text{ch}(j).$$

The vertices j in step 2 can be ordered in any order. However, by defining $V_j = S_{I_{jj}}$ as in Algorithm 4.1 and using a reverse topological ordering, we can avoid having to extract $S_{I_{jj}}$ from the CCS structure of S . In the modified algorithm, the second step is replaced by

$$M_{jj} = \frac{K_{jj}}{D_{jj}^2}, \quad M_{I_{jj}} = \frac{1}{D_{jj}} V_j K_{I_{jj}}, \quad V_i = E_{J_j I_i}^T \begin{bmatrix} S_{jj} & S_{I_{jj}}^T \\ S_{I_{jj}} & V_j \end{bmatrix} E_{J_j I_i}, \quad i \in \text{ch}(j),$$

for $j \in \{1, 2, \dots, n\}$ in a reverse topological order.

5.3 Inverse Hessian

To evaluate $Y = \mathcal{H}^{-1}(T)$, we use the equation (30) to compute the linearized Cholesky factors D' , L' from T , and the equation (31) to compute Y from D' , L' . We use the same notation (32) as in the previous section.

Algorithm 5.2. *Inverse Hessian evaluation.*

Input. A matrix $T \in \mathbf{S}_V^n$, the Cholesky L , D of a positive definite matrix $X \in \mathbf{S}_V^n$, and the projected inverse $S = \mathcal{P}(X^{-1})$.

Output. The matrix $Y = \mathcal{H}^{-1}(T)$, i.e., the solution $Y \in \mathbf{S}_V^n$ of the equation $\mathcal{H}(Y) = \mathcal{P}(X^{-1}YX^{-1}) = T$.

Algorithm.

1. Iterate over $j \in \{1, 2, \dots, n\}$ in reverse topological order. For each j , calculate the j th column of M from

$$\begin{bmatrix} M_{jj} & M_{I_{jj}}^T \\ M_{I_{jj}} & V_j' \end{bmatrix} = \begin{bmatrix} 1 & L_{I_{jj}}^T \\ 0 & I \end{bmatrix} \begin{bmatrix} T_{jj} & T_{I_{jj}}^T \\ T_{I_{jj}} & V_j' \end{bmatrix} \begin{bmatrix} 1 & 0 \\ L_{I_{jj}} & I \end{bmatrix}$$

and the update matrices V_i' for the children of vertex j via

$$V_i' = E_{J_j I_i}^T \begin{bmatrix} T_{jj} & T_{I_{jj}}^T \\ T_{I_{jj}} & V_j' \end{bmatrix} E_{J_j I_i}, \quad i \in \text{ch}(j).$$

2. For $j \in \{1, 2, \dots, n\}$, compute column j of K via

$$K_{jj} = D_{jj}^2 M_{jj}, \quad K_{I_{jj}} = D_{jj} S_{I_j I_j}^{-1} M_{I_{jj}}.$$

3. Iterate over $j \in \{1, 2, \dots, n\}$ in topological order. For each j , compute U_j' and the j th column of Y from

$$\begin{bmatrix} Y_{jj} & Y_{I_{jj}}^T \\ Y_{I_{jj}} & -U_j' \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ L_{I_{jj}} & I \end{bmatrix} \begin{bmatrix} K_{jj} & K_{I_{jj}}^T \\ K_{I_{jj}} & 0 \end{bmatrix} \begin{bmatrix} 1 & L_{I_{jj}}^T \\ 0 & I \end{bmatrix} - \sum_{i \in \text{ch}(j)} E_{J_j I_i} U_i' E_{J_j I_i}^T.$$

A improvement of step 2 is to use a factorization $S_{I_j I_j} = -R_j R_j^T$ and update the matrices R_j recursively, following a reverse topological order, as discussed at the end of Section 4.

5.4 Hessian factor

Step 1 in the Hessian evaluation algorithm (Algorithm 5.1) is a linear mapping that transforms Y to K . Step 3 is a linear mapping that transforms M to T . It is interesting to note that these two mappings are adjoints. Step 2 implements a self-adjoint and positive definite mapping, which transforms K to M . Factoring the positive definite mapping in step 2 provides a factorization

$$\mathcal{H}(Y) = \mathcal{R}^{\text{adj}}(\mathcal{R}(Y)),$$

with \mathcal{R} a linear mapping from \mathbf{S}_V^n to \mathbf{S}_V^n . The factorization of the mapping in step 2 can be implemented by defining a factorization $S_{I_j I_j} = R_j R_j^T$ with R_j upper triangular for each vertex j . Although it is impractical to pre-compute and store the matrices R_j for each vertex, they can be efficiently computed recursively in a reverse topological order, as in Algorithm 4.3. For the sake of clarity, we omit the details in the following algorithms.

The algorithm for evaluating \mathcal{R} consists of step 1 in Algorithm 5.1 and one half of step 2.

Algorithm 5.3. *Evaluation of Hessian factor.*

Input. A matrix $Y \in \mathbf{S}_V^n$, the Cholesky factors L, D of a positive definite matrix $X \in \mathbf{S}_V^n$, and the projected inverse $S = \mathcal{P}(X^{-1})$.

Output. The matrix $W = \mathcal{R}(Y)$ where $\mathcal{H} = \mathcal{R}^{\text{adj}} \circ \mathcal{R}$ is the Hessian of f at X .

Algorithm.

1. Iterate over $j \in \{1, 2, \dots, n\}$ in topological order. For each j , calculate the j th column of K , and the update matrix U'_j via

$$\begin{bmatrix} K_{jj} & K_{I_j j}^T \\ K_{I_j j} & U'_j \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ -L_{I_j j} & I \end{bmatrix} \left(\begin{bmatrix} Y_{jj} & Y_{I_j j}^T \\ Y_{I_j j} & 0 \end{bmatrix} + \sum_{i \in \text{ch}(j)} E_{J_j I_i} U'_i E_{J_j I_i}^T \right) \begin{bmatrix} 1 & -L_{I_j j}^T \\ 0 & I \end{bmatrix}.$$

2. For all $j \in \{1, 2, \dots, n\}$, compute column j of W via

$$W_{jj} = \frac{K_{jj}}{D_{jj}}, \quad W_{I_j j} = \frac{1}{\sqrt{D_{jj}}} R_j^T K_{I_j j}$$

where R_j is a triangular factor of $S_{I_j I_j} = R_j R_j^T$.

The algorithm for evaluating \mathcal{R}^{adj} consists of the second half of step 2 of Algorithm 5.1 and of step 3. It also readily follows by taking the adjoint of the calculations in Algorithm 5.3.

Algorithm 5.4. *Evaluation of adjoint Hessian factor.*

Input. A matrix $W \in \mathbf{S}_V^n$, the Cholesky factors L, D of a positive definite $X \in \mathbf{S}_V^n$, and the projected inverse $S = \mathcal{P}(X^{-1})$.

Output. The matrix $T = \mathcal{R}^{\text{adj}}(W)$ where $\mathcal{H} = \mathcal{R}^{\text{adj}} \circ \mathcal{R}$ is the Hessian of f at X .

Algorithm.

1. For all $j \in \{1, 2, \dots, n\}$, compute column j of M via

$$M_{jj} = \frac{W_{jj}}{D_{jj}}, \quad M_{I_j j} = \frac{1}{\sqrt{D_{jj}}} R_j W_{I_j j}.$$

2. Iterate over $j \in \{1, 2, \dots, n\}$ in reverse topological order. For each j , calculate the j th column of T from

$$\begin{bmatrix} T_{jj} & T_{I_j j}^T \\ T_{I_j j} & V'_j \end{bmatrix} = \begin{bmatrix} 1 & -L_{I_j j}^T \\ 0 & I \end{bmatrix} \begin{bmatrix} M_{jj} & M_{I_j j}^T \\ M_{I_j j} & V'_j \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -L_{I_j j} & I \end{bmatrix}$$

and the update matrices V'_i for the children of vertex j via

$$V'_i = E_{J_j I_i}^T \begin{bmatrix} T_{jj} & T_{I_j j}^T \\ T_{I_j j} & V'_j \end{bmatrix} E_{J_j I_i}, \quad i \in \text{ch}(j).$$

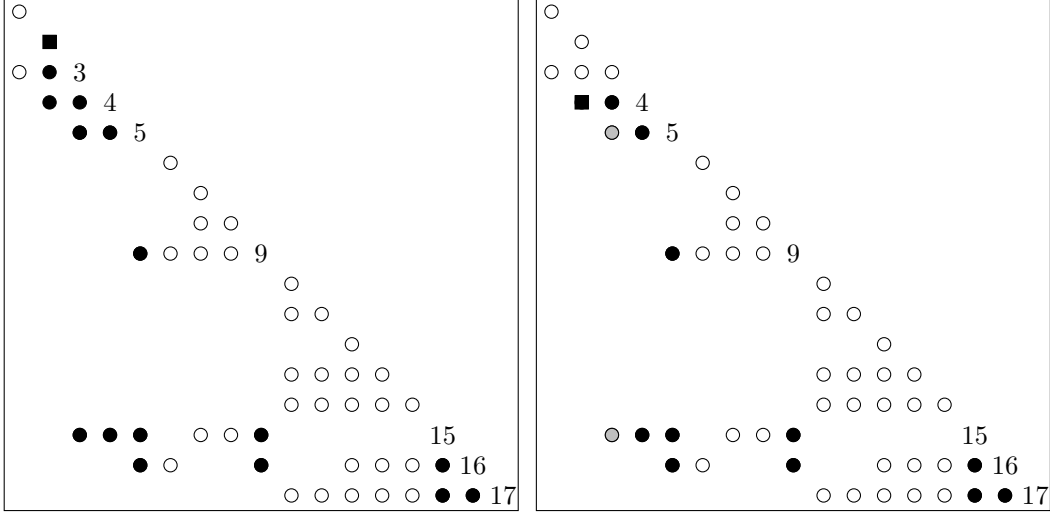


Figure 7: Sparsity pattern of the lower triangular part of $W = \mathcal{R}(Y)$ when (a) Y has a single diagonal nonzero in position (2,2), marked with a solid square, and (b) when Y has an off-diagonal nonzero in position (4,2). The solid entries mark the nonzero elements in W . The gray markers in (b) correspond to nonzeros in W that are introduced by the scaling in step 2 of Algorithm 5.3, *i.e.*, these entries are not present in the intermediate variable K .

5.5 Sparse arguments

In many applications, such as interior-point methods for the cone programs (7) mentioned in the introduction, the Hessian $\mathcal{H} = \nabla^2 f(X)$ and its factorization $\mathcal{H} = \mathcal{R}^{\text{adj}} \circ \mathcal{R}$ are needed to compute coefficients

$$H_{ij} = A_i \bullet (\nabla f(X)[A_j]) = \mathcal{R}(A_i) \bullet \mathcal{R}(A_j)$$

for m matrices $A_1, \dots, A_m \in \mathbf{S}_V^n$. The matrices A_i are often very sparse relative to the sparsity pattern V . In this section, we examine the implications of sparsity in the matrix Y on the computation of $W = \mathcal{R}(Y)$ using Algorithm 5.3.

From step 1 in Algorithm 5.3 we see that if $Y_{jj} \neq 0$, then $K_{I_{jj}}$ and U'_j are nonzero and dense. In step 1 these nonzeros are then further propagated via the recursion in topological order to all the columns indexed by ancestors of j . Therefore, if $Y_{jj} \neq 0$, then $W_{J_k k} \neq 0$ are dense for all ancestors k of j .

To see how an off-diagonal nonzero in Y affects the sparsity pattern of W , suppose that $Y_{jj} = 0$, $Y_{ij} \neq 0$ for some $i \in I_j$, and $U'_k = 0$ for all $k \in \text{ch}(j)$. Then $W_{pq} \neq 0$ for all ancestors q of j and $p \in J_j \cap \{i, i+1, \dots, n\}$. Hence nonzeros in column j of Y create fill in the columns that correspond to ancestors of j . In Algorithm 5.3, we can therefore prune the elimination tree at node k if all the descendants of node k correspond to columns in Y with no lower triangular nonzeros.

Figure 7 shows examples of the sparsity pattern of W when Y has a diagonal and an off-diagonal nonzero, respectively, for the pattern V in Figure 1.

Numerical results To evaluate the benefits of exploiting additional sparsity in the argument, we have implemented and tested a version of Algorithm 5.3 that exploits sparsity in Y relative to V . In the experiment, we use as test data a set of randomly generated problems with sparsity

Sparsity pattern	n	Dense (seconds)	Sparse (seconds)	Ratio
HB/plat1919	1919	0.17	0.05	3.7
HB/bcsstk13	2003	1.40	0.36	3.9
HB/lshp3025	3025	0.21	0.05	4.6
Boing/nasa4704	4704	1.05	0.27	4.0
TKK/g3rmt3m3	5357	1.43	0.35	4.1
Schenk_IBMNA/c-36	7479	0.37	0.06	6.1
Wang/swang1	10800	7.05	1.72	4.1
ACUSIM/Pres_Poisson	14822	17.65	4.35	4.1
GHS_psdef/wathen100	30401	6.05	1.38	4.4

Table 1: Average computational times (10 trials) for Algorithm 5.3 with and without the technique for sparse arguments (columns ‘sparse’ and ‘dense’, respectively) when applied to sparse arguments. The sparse arguments are randomly generated with two lower-triangular nonzero entries in random positions.

patterns from the University of Florida Sparse Matrix Collection. For each sparsity pattern, we generate ten sparse arguments with just two lower-triangular nonzero entries in random positions. Table 1 summarizes the average time required to compute $W = \mathcal{R}(Y)$ using Algorithm 5.3 with and without the techniques described in this section. For these very sparse arguments, pruning the elimination tree results in average speedups in the range 4–6, but in general the speedup depends both on the number of nonzeros and on the position of the nonzeros.

Finally, we remark that additional computational savings can be made in step 2 of Algorithm 5.3 when $\mathcal{R}(Y)$ is needed for several arguments Y . Specifically, the triangular factors R_j of $S_{I_j I_j} = R_j R_j^T$ need only be computed once.

6 Cliques and clique trees

It is known that the performance of sparse Cholesky factorization algorithms on modern computers can be improved by combining groups of vertices into *supernodes* and applying block elimination to the corresponding columns. Several definitions of supernodes exist in the literature. In this paper, we define a supernode as a maximal group of columns of L (sorted, but not necessarily contiguous) that share the same nonzero structure. More specifically, if N is a supernode and $j = \max N$, then for all $k \in N$,

$$J_k = (N \cup I_j) \cap \{k, k+1, \dots, n\}.$$

In the example of Figure 1, the sets

$$\{1\}, \quad \{2\}, \quad \{3, 4\}, \quad \{5, 9\}, \quad \{6\}, \quad \{7, 8\}, \quad \{10, 11\}, \quad \{12, 13, 14\}, \quad \{15, 16, 17\}$$

form supernodes. (Another more common definition adds the requirement that the indices in N are contiguous [LNP93]; this can be achieved from a set of supernodes as defined above by a simple reordering.) In the context of multifrontal factorizations, the grouping into supernodes has the advantage that only one frontal matrix is required per supernode. This reduces the memory and arithmetic overhead incurred for the assembly of frontal matrices. Moreover, the block operations

allow us to replace matrix-vector operations (level-2 BLAS) with more efficient matrix-matrix operations (level-3 BLAS) [DCHD90].

Supernodes are closely related to cliques in the filled graph and the barrier algorithms described in [DVR08, DV09], which involve iterations on clique trees, can be interpreted as supernodal multifrontal algorithms. We therefore start the discussion with a review of cliques and clique trees, and their connections with supernodes.

6.1 Cliques

A filled graph is also known as a *chordal* or *triangulated* graph. A *clique* is a maximal set of vertices that define a complete subgraph of the filled graph. Equivalently, a clique is a set of indices that define a dense lower-triangular principal subblock of L . Every clique W in a filled graph can be expressed as $W = J_i$, where $i = \min W$, the least element in W [BP93, proposition 2]. This follows from the fact that the index sets J_i define complete subgraphs, as noted in Section 2. Hence, if i is the lowest index in the clique W , then $W \subseteq J_i$. Since W is maximal, we must have $W = J_i$. The vertex $i = \min W$ is called the *representative vertex* of the clique. Since there are at most n representative vertices, a filled graph can have at most n cliques.

Efficient algorithms for identifying the representative vertices can be derived from the following criterion: the cliques are exactly the sets J_i for which there exists no $j < i$ with $J_i \subset J_j$ [LPP89, proposition 3]. This follows from the characterization of cliques in terms of representative vertices. If $J_i \subset J_j$ for some $j < i$, then J_i is certainly not a clique, since it is strictly included in another complete subgraph. Conversely, if J_i is not a clique, *i.e.*, $J_i \subset W$ for some clique W , and j is the representative vertex of W , then $j < i$ and $J_i \subset W = J_j$.

In the example in Figure 1, the representative vertices are 1, 2, 3, 5, 6, 7, 10, 12, 15. The other vertices are not representative because the corresponding sets J_j are not maximal:

$$J_4 \subset J_3, \quad J_8 \subset J_7, \quad J_9 \subset J_5, \quad J_{11} \subset J_{10}, \quad J_{13}, J_{14}, J_{16}, J_{17} \subset J_{12}.$$

The representative vertices are easily identified from the elimination tree and the monotone degrees of the vertices. It can be shown that a vertex j is a representative vertex if and only if

$$|I_j| > |I_k| - 1 \quad \forall k \in \text{ch}(j) \quad (33)$$

(see [PS90]). To see this, recall from (9) that $I_k \subseteq J_j$ and $|I_j| \geq |I_k| - 1$ if $k \in \text{ch}(j)$. Therefore, if $|I_j| = |I_k| - 1$ for some $k \in \text{ch}(j)$, then $J_j = I_k \subset J_k$. Therefore j is not a representative vertex because the complete subgraph defined by J_j is not maximal. Conversely, suppose j is not representative, *i.e.*, $J_j \subset J_l$ for some $l < j$. In particular, $L_{jl} \neq 0$ and therefore, l is a descendant of j in the elimination tree. From Theorem 2 (the set of inequalities (10)), this implies that $J_j \subset J_k$ for all k in the path from l to j . In particular,

$$J_j \subset J_k \subseteq J_j \cup \{k\}$$

for the child k of j on this path. Therefore, $J_j = I_k$ and $|I_j| = |I_k| - 1$.

6.2 Supernode partitions

By comparing the monotone degrees of the vertices in the elimination tree and their parents we can partition the vertices $\{1, 2, \dots, n\}$ in sets $N_j = \{i_1, i_2, \dots, i_r\}$ where $i_1 = j$ is a representative

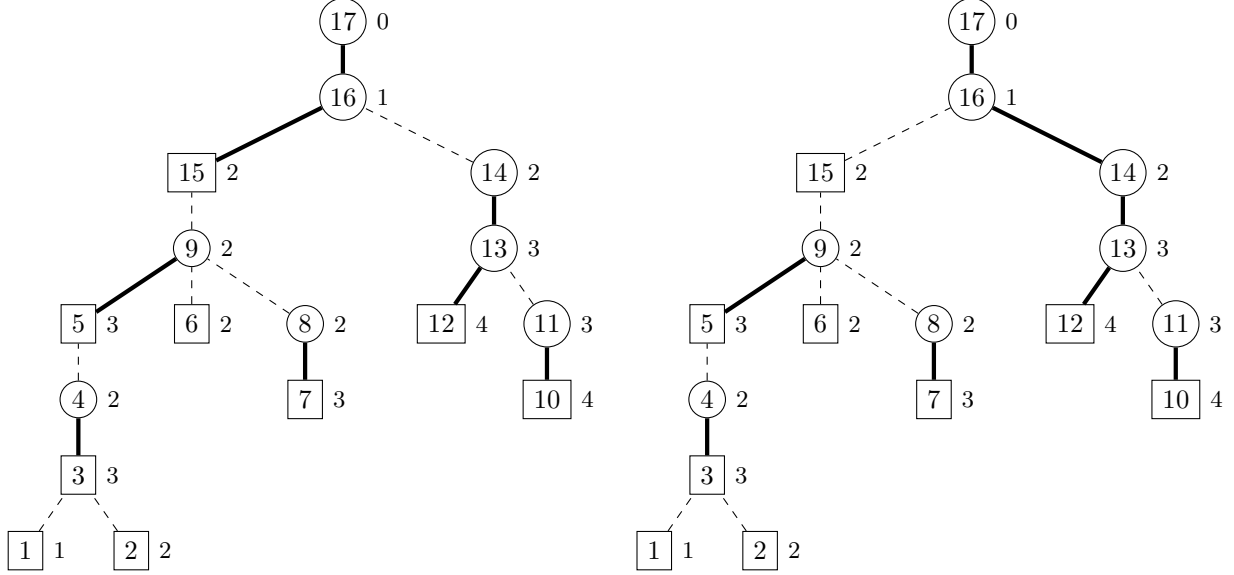


Figure 8: Two supernode partitions of the elimination tree in Figure 1. The representative vertices are shown as rectangles. The vertices joined by the paths shown with heavy lines form the sets N_k , where k is the representative vertex on the path. The sets N_k are enumerated in Table 2.

clique vertex. The vertices i_1, i_2, \dots, i_r form a path from j to an ancestor i_r of j in the elimination tree, and

$$|I_{i_1}| = |I_{i_2}| + 1 = |I_{i_3}| + 2 = \dots = |I_{i_r}| + r - 1,$$

or, equivalently,

$$I_{i_1} = J_{i_2} = \{i_2\} \cup J_{i_3} = \dots = \{i_2, i_3, \dots, i_{r-1}\} \cup J_{i_r}. \quad (34)$$

The sets N_j are supernodes (in the definition given at the beginning of this section). In general, several such partitions exist. Two possible partitions for the elimination tree in Figure 1 are shown in Figure 8 and listed in Table 2. The representative vertices are shown as rectangles, and the sets N_j are the vertices on the paths shown with heavy lines. We note two important properties of the sets N_j :

- The set N_j is a subset of the clique represented by vertex j : $N_j \subseteq J_j$. This can be seen from (34) which implies $J_i \subset J_j$ if $i \in N_j$ and $i \neq j$.
- Define $A_j = J_j \setminus N_j$. Then we have $i < \min A_j$ for all $i \in N_j$. To see this, first note that if $k \in J_j$, then $L_{kj} \neq 0$ and therefore, k is an ancestor of j in the elimination tree (Theorem 1). If also $k \leq i$ for some $i \in N_j$, then k is on the path from j to i in the elimination tree. However, by definition of N_j , this means that $k \in N_j$.

This result means that $N_j \cup A_j$ is an ordered partition of the clique J_j , *i.e.*, the elements of N_j have a lower index than the elements of A_j . (In [LPP89, PS90], the sets N_j and A_j are referred to as the *new set* $\text{new}(K)$ and the *ancestor set* $\text{anc}(K)$, respectively, where K is the clique $K = J_j$.)

If A_j is nonempty, we refer to the vertex $k = \min A_j$ as the *first ancestor* of the clique J_j . The first ancestor can be identified from the elimination tree as the parent of the vertex $i = \max N_j$. This follows from (34) with $i = i_r$ and the fact that the parent of vertex i is the first element in

k	N_k	A_k	k	N_k	A_k
1	{1}	{3}	1	{1}	{3}
2	{2}	{3, 4}	2	{2}	{3, 4}
3	{3, 4}	{5, 15}	3	{3, 4}	{5, 15}
5	{5, 9}	{15, 16}	5	{5, 9}	{15, 16}
6	{6}	{9, 16}	6	{6}	{9, 16}
7	{7, 8}	{9, 15}	7	{7, 8}	{9, 15}
10	{10, 11}	{13, 14, 17}	10	{10, 11}	{13, 14, 17}
12	{12, 13, 14}	{16, 17}	12	{12, 13, 14, 15, 16, 17}	{}
15	{15, 16, 17}	{}	15	{15}	{16, 17}

Table 2: The two supernode partitions defined in Figure 8. The first columns of the tables are the representative vertices. Each clique J_k is partitioned in two sets as $J_k = N_k \cup A_k$. The sets N_k for a partition of the vertices $\{1, 2, \dots, n\}$

J_i greater than i . Note that while the first ancestor can be determined from the elimination tree and the sets N_j , the rest of the sets A_j cannot be derived from the elimination tree but require knowledge of the clique J_j .

The sets N_k and A_k for the two partitions in the example are listed in Table 2.

6.3 Clique trees

We can associate with the vertex partitioning in sets N_j a tree with the cliques J_j as its nodes. The root of the clique tree is the clique represented by the vertex j for which $n \in N_j$. The parent of the clique J_j is the clique which has as its representative the vertex k for which $\min A_j \in N_k$. We will use the notation $k = \text{par}(j)$ to denote that J_k is the parent of J_j in the clique tree. Figure 9 shows the clique trees defined by the partitions N_j in Figure 8. The clique tree satisfies the following key properties [PS90, p.186] [LPP89]:

- $A_j \subset J_{\text{par}(j)}$.

Indeed, let $i = \min A_j$ be the first ancestor of clique J_j . By definition of the clique tree, $i \in N_{\text{par}(j)}$. From the definition in (34) this implies that $J_i \subset J_{\text{par}(j)}$. Since A_j defines a complete subgraph of vertices with indices greater than or equal to i , we have $A_j \subseteq J_i$. Therefore $A_j \subset J_{\text{par}(j)}$.

- An element of N_j is in the clique J_k only if J_k is a descendant of J_j in the clique tree.

We can show this by contradiction. Suppose $i \in N_j$ belongs to J_k and J_k is not a descendant of J_j in the clique tree. The sets N_l form a partition of $\{1, 2, \dots, n\}$, so if $i \in N_j$ and $i \in J_k$ for $k \neq j$, then $i \in A_k$. By the previous property, this implies $i \in J_{\text{par}(k)}$. We have $\text{par}(k) \neq j$ because J_k is not a descendant of J_j . Therefore $i \in A_{\text{par}(k)}$ and, again from the previous property, $i \in J_{\text{par}(\text{par}(k))}$ and $i \in A_{\text{par}(\text{par}(k))}$. Continuing this process recursively, we eventually arrive at the conclusion that i belongs to A_r where J_r is the root of the clique tree. However, this is impossible because $J_r = N_r$ and $A_r = \emptyset$.

- If an element of N_j is in A_k , then it belongs to all the cliques on the path between J_j and J_k in the clique tree.

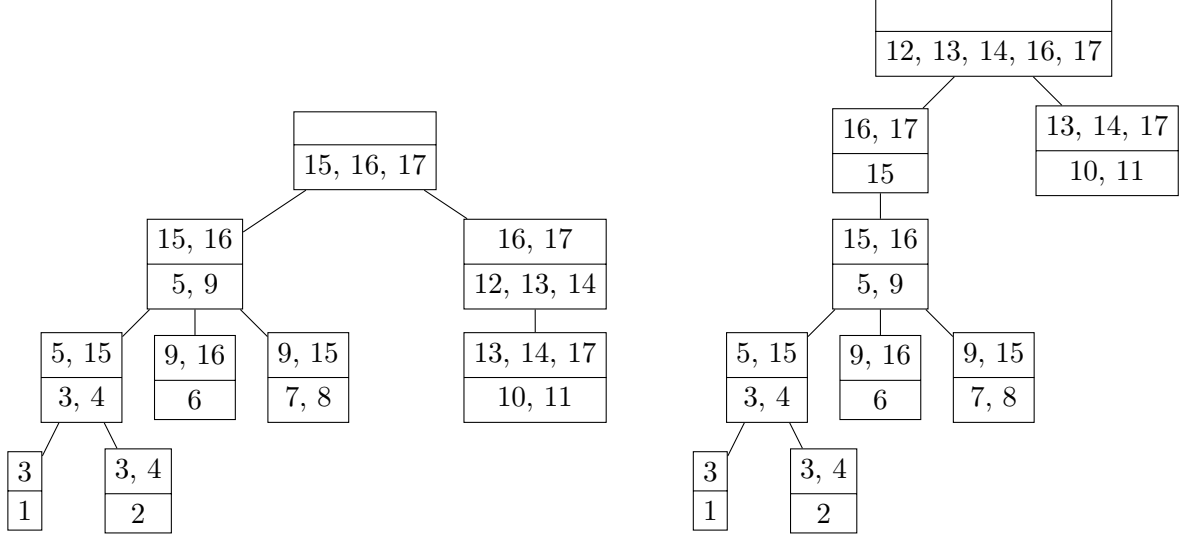


Figure 9: The two clique trees corresponding to the partitions in Figure 8. The first index in the bottom row at each node is the representative vertex j of the clique J_j . Each clique J_j is partitioned in two sets $J_j = N_j \cup A_j$. The indices of the bottom row at each node form the sets N_j ; the indices in the top row form the set A_j . The parent of clique J_j is the clique J_k that includes the first ancestor of clique J_j in its N_k set.

This follows by combining the first two properties. From the second property, if $i \in N_j$ and $i \in A_k$, then J_k is a descendant of J_j . Assume there are cliques $J_s, J_{\text{par}(s)}$ on the path between J_k and J_j with the property that $i \in J_s$ and $i \notin J_{\text{par}(s)}$. From the first property, this implies $i \in N_s$. But this contradicts $i \in N_j$, unless $j = s$, because $N_j \cup N_s = \emptyset$ if $j \neq s$.

Taken together, these three properties state that the cliques that contain a vertex i form a subtree in the clique tree. The root of the subtree is the unique clique J_k for which $i \in N_k$. This is known as the *induced subtree property* of clique trees [BP93].

6.4 Clique tree algorithm

To summarize the results of this section, we state a simple algorithm that identifies the representative vertices of the cliques, generates a partition into sets N_k , identifies the first ancestors $\min A_k$ of the cliques, and determines the parent structure of the clique tree. The algorithm is due to Pothén and Sun [PS90, p.185].

Algorithm 6.1. Clique tree algorithm.

Input. An elimination tree and the monotone degree $|I_k|$, $k = 1, \dots, n$.

Output. The representative vertices, the partition in supernodes N_j , the first ancestor of each clique, and the parent structure of a clique tree.

Algorithm. For $i = 1, \dots, n$:

1. If $|I_i| > |I_j| - 1$ for all $j \in \text{ch}(i)$, then i is a representative vertex. Set $N_i = \{i\}$ and $k = i$. Otherwise, choose a vertex $j \in \text{ch}(i)$ with $|I_i| = |I_j| - 1$, determine the representative vertex k for which $j \in N_k$, and add i to N_k .
2. For each $j \in \text{ch}(i)$, if $j \in N_l$ and $l \neq k$, set $\text{par}(l) := k$ and $\min A_l := i$.

Note that in step 1, there may be several choices for the child vertex j , and these choices lead to different vertex partitions and different clique trees. The vertex $i = 16$ in Figure 8, for example, has two children j that both satisfy $|I_i| = |I_j| - 1$. These two choices lead to the different vertex partitions in Figure 8 and the two clique trees in Figure 9.

7 Supernodal multifrontal algorithms

We assume there are l cliques, with representative nodes i_1, \dots, i_l and that the sets N_i are contiguous. The supernodal algorithms are block versions of the multifrontal algorithms in which the scalar diagonal elements X_{jj} are replaced with dense principal blocks $X_{N_j N_j}$ and the subcolumns $X_{I_j j}$ with dense submatrices $X_{A_j N_j}$.

For a clique J_k we denote by $\text{ch}(J_k)$ the set of child cliques of J_k in the clique tree. This is not to be confused with $\text{ch}(k)$ (with a vertex k as argument), which refers to the children of the vertex k in the elimination tree.

In this section we start with a supernodal version of the Cholesky factorization algorithm. We then give similar extensions of the primal and dual gradient evaluation algorithms. For the sake of brevity, we will omit the extensions of the other algorithms in Sections 3–5, which follow the same pattern.

7.1 Cholesky factorization

In the supernodal Cholesky factorization we factor X as $X = LDL^T$ with D block-diagonal and L unit lower triangular. The matrix D has l dense diagonal blocks $D_{N_i N_i}$ for $i \in \{i_1, \dots, i_l\}$. Corresponding with each clique, L has a diagonal block $L_{N_i N_i} = I$ and a dense submatrix $L_{A_i N_i}$. The rest of the block-column indexed by N_i is zero.

Algorithm 7.1. Cholesky factorization.

Input. A positive definite matrix $X \in \mathbf{S}_V^n$ and a clique tree for the sparsity pattern V .

Output. The factors L, D in the Cholesky factorization $X = LDL^T$.

Algorithm. Iterate over $j \in \{i_1, i_2, \dots, i_l\}$ using a topological order of the clique tree.

For each j , form the frontal matrix

$$F_j = \begin{bmatrix} F_{11} & F_{21}^T \\ F_{21} & F_{22} \end{bmatrix} = \begin{bmatrix} X_{N_j N_j} & X_{A_j N_j}^T \\ X_{A_j N_j} & 0 \end{bmatrix} + \sum_{J_i \in \text{ch}(J_j)} E_{J_j A_i} U_i E_{J_j A_i}^T$$

and calculate $D_{N_j N_j}$, $L_{A_j N_j}$, and the update matrix U_j from

$$D_{N_j N_j} = F_{11}, \quad L_{A_j N_j} = F_{21} D_{N_j N_j}^{-1}, \quad U_j = F_{22} - L_{A_j N_j} D_{N_j N_j} L_{A_j N_j}^T.$$

As can be seen, only one frontal matrix is assembled per clique, a major advantage compared to Algorithm 3.1. Moreover, the main computation is the level-3 BLAS operation in the computation of U_j .

7.2 Gradients

The supernodal counterpart of Algorithms 4.1 for computing the primal gradient or projected inverse is as follows. In this algorithm, V_j is a dense ‘update matrix’ defined as $V_j = S_{A_j A_j}$.

Algorithm 7.2. Projected inverse.

Input. The Cholesky factors L, D of a positive definite matrix $X \in \mathbf{S}_V^n$.

Output. The projected inverse $S = \mathcal{P}(X^{-1}) = -\nabla f(X)$.

Algorithm. Iterate over $j \in \{i_1, \dots, i_l\}$ using a reverse topological order of the clique tree. For each j , calculate $S_{N_j N_j}$ and $S_{A_j N_j}$ from

$$S_{A_j N_j} = -V_j L_{A_j N_j}, \quad S_{N_j N_j} = D_{N_j N_j}^{-1} - S_{A_j N_j}^T L_{A_j N_j} \quad (35)$$

and compute the update matrices

$$V_i = E_{J_j A_i}^T \begin{bmatrix} S_{N_j N_j} & S_{A_j N_j}^T \\ S_{A_j N_j} & V_j \end{bmatrix} E_{J_j A_i}, \quad J_i \in \text{ch}(J_j). \quad (36)$$

The main calculation is the matrix-matrix product in (35) which replaces the matrix-vector product (23) in the multifrontal algorithm.

The extension of Algorithm 4.2 for computing the dual gradient or the maximum determinant positive definite completion is as follows.

Algorithm 7.3. Matrix completion.

Input. A matrix $S \in \mathbf{S}_V^n$ that has a positive definite completion.

Output. The Cholesky factors L, D of $X = -\nabla f_*(S)$, *i.e.*, of the positive definite matrix $X \in \mathbf{S}_V^n$ that satisfies $\mathcal{P}(X^{-1}) = S$.

Algorithm. Iterate over $j \in \{i_1, \dots, i_l\}$ using a reverse topological order of the clique tree. For each j , compute $D_{N_j N_j}$ and $L_{A_j N_j}$ from

$$L_{A_j N_j} = -V_j^{-1} S_{A_j N_j}, \quad D_{N_j N_j} = (S_{N_j N_j} + S_{A_j N_j}^T L_{A_j N_j})^{-1}, \quad (37)$$

and compute the update matrices

$$V_i = E_{J_j A_i}^T \begin{bmatrix} S_{N_j N_j} & S_{A_j N_j}^T \\ S_{A_j N_j} & V_j \end{bmatrix} E_{J_j A_i}, \quad J_i \in \text{ch}(J_j). \quad (38)$$

As for the multifrontal completion algorithm with factored update matrices (Algorithm 4.3), this algorithm can be improved by propagating a factorization of V_j and using (38) to compute the factors of V_i from the factors of V_j . We mentioned in section 6.2 that for every clique, the vertices in N_i precede those in A_i . As a consequence, the matrix $E_{J_j A_i}$ can be partitioned as

$$E_{J_j A_i} = \begin{bmatrix} E_{N_j, A_i \cap N_j} & 0 \\ 0 & E_{A_j, A_i \setminus N_j} \end{bmatrix}.$$

Using this property in (38) we get

$$V_i = E_{J_j A_i}^T \begin{bmatrix} S_{N_j N_j} & S_{A_j N_j}^T \\ S_{A_j N_j} & V_j \end{bmatrix} E_{J_j A_i} = \begin{bmatrix} A & B^T \\ B & CC^T \end{bmatrix}$$

where A is a principal submatrix of $S_{N_j N_j}$ of order $|A_i \cap N_j|$, and C consists of $|A_i \setminus N_j|$ rows of the upper triangular factor R_j of $V_j = R_j R_j^T$. By reducing C to square triangular form $C = RQ^T$ using a series of Householder transformations, and a factorization $A - B^T R^{-T} R^{-1} B = \tilde{R} \tilde{R}^T$ we obtain the factorization $V_i = R_i R_i^T$ as

$$R_i = \begin{bmatrix} \tilde{R} & (R^{-1} B)^T \\ 0 & R \end{bmatrix}.$$

7.3 Numerical results

We apply the supernodal multifrontal algorithms to the test problems described in Section 4.3. The left-hand plot in Figure 10 shows the CPU times for Algorithms 7.1 (Cholesky factorization) and 7.2 (projected inverse or primal gradient). The right-hand plot shows the CPU times for Algorithms 7.2 and 4.1 (supernodal multifrontal and multifrontal projected inverse, respectively). From the first plot we see that the cost of evaluating the primal gradient is comparable to the cost of computing the Cholesky factorization. The second plot shows that the supernodal implementation of the primal gradient is substantially faster than the non-supernodal implementation, for all but a few of the small problems.

8 Conclusions

We have derived recursive algorithms for evaluating the values, gradients, and Hessians of the primal and dual barriers

$$f(X) = -\log \det X, \quad f_*(S) = \sup_{X \in \mathbf{S}_V^n} (-\mathbf{tr}(SX) - f(X)),$$

defined for sparse symmetric matrices $X, S \in \mathbf{S}_V^n$ with a given sparsity pattern V , where V is a filled (or chordal) pattern. Our interest in these algorithms is motivated by their importance in interior-point methods for conic optimization with sparse matrix cone constraints [ADV10]. Similar algorithms can be formulated for closely related problems that arise in sparse semidefinite programming, for example, the matrix completion techniques used in primal-dual methods [FKMN00, NFF⁺03].

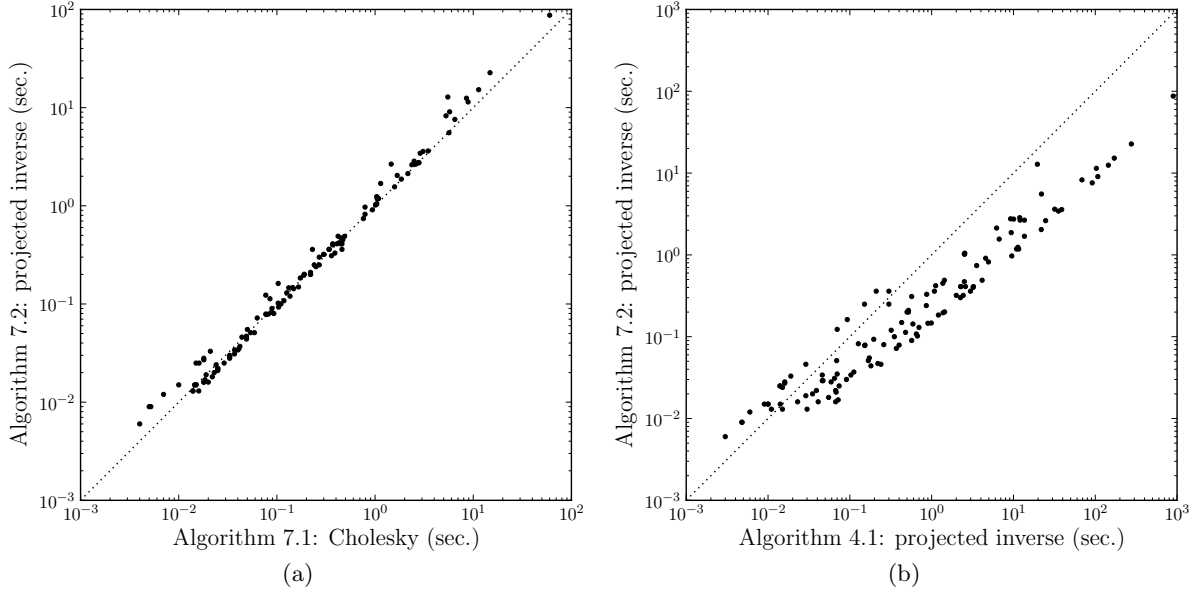


Figure 10: Scatter plots of CPU times for supernodal algorithms applied to 128 test problems. The plot on the left shows that the cost of computing the Cholesky factorization and the projected inverse is approximately the same. The plot on the right shows that the supernodal multifrontal implementation of the projected inverse algorithm is faster than the multifrontal implementation for all but some small problems.

Our goal was to formulate efficient barrier algorithms based on Cholesky factorization techniques for large sparse matrices and, specifically, the multifrontal algorithm that has been extensively studied in the sparse matrix literature since the 1980s. The algorithms inherit many of the properties of the multifrontal method. This means that a wide range of known techniques from the sparse matrix literature can be used to further improve the algorithms. For example, tree parallelism and node parallelism are readily exploited in a multifrontal method [ADL00]. Relaxed supernodes and supernode amalgamation techniques [DR83, AG89] have also been shown to improve the performance. Other improvements include tree modifications [Liu88] and memory optimization techniques [Liu86, GL06].

The starting point in this paper was the multifrontal Cholesky factorization algorithm. Similar algorithms can be derived from the other popular types of sparse factorization algorithms, such as the up-looking Cholesky factorization (used in CHOLMOD [CDHR08]) or the left-looking Cholesky factorization (a blocked version of which is used in CHOLMOD’s supernodal solver). It would be of interest to compare the performance of these algorithms with the multifrontal algorithms formulated in this paper.

Acknowledgment

This material is based upon work supported by the National Science Foundation under Grants No. ECCS-0824003 and DMS-1115963. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the

National Science Foundation.

References

- [ADL00] P. R. Amestoy, I. S. Duff, and J.Y. L'Excellent. Multifrontal parallel distributed symmetric and unsymmetric solvers. *Computer Methods in Applied Mechanics and Engineering*, 184:501–520, 2000.
- [ADR⁺10] P. R. Amestoy, I. S. Duff, Y. Robert, F.-H. Rouet, and B. Uçar. On computing inverse entries of a sparse matrix in an out-of-core environment. Technical Report TR/PA/10/59, CERFACS, 2010.
- [ADV10] M. S. Andersen, J. Dahl, and L. Vandenbergh. Implementation of nonsymmetric interior-point methods for linear optimization over sparse matrix cones. *Mathematical Programming Computation*, 2:167–201, 2010.
- [AG89] C. Ashcraft and R. Grimes. The influence of relaxed supernode partitions on the multifrontal method. *ACM Transactions on Mathematical Software*, 15(4):291–309, 1989.
- [BEd08] O. Banerjee, L. El Ghaoui, and A. d'Aspremont. Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data. *Journal of Machine Learning Research*, 9:485–516, 2008.
- [BP93] J. R. S. Blair and B. Peyton. An introduction to chordal graphs and clique trees. In A. George, J. R. Gilbert, and J. W. H. Liu, editors, *Graph Theory and Sparse Matrix Computation*. Springer-Verlag, 1993.
- [Bur03] S. Burer. Semidefinite programming in the space of partial positive semidefinite matrices. *SIAM Journal on Optimization*, 14(1):139–172, 2003.
- [CD95] Y. E. Campbell and T. A. Davis. Computing the sparse inverse subset: an inverse multifrontal approach. Technical Report TR-95-021, Computer and Information Sciences Department, University of Florida, 1995.
- [CDHR08] Y. Chen, T. A. Davis, W. W. Hager, and S. Rajamanickam. Algorithm 887: CHOLMOD, supernodal sparse Cholesky factorization and update/downdate. *ACM Transactions on Mathematical Software*, 35(3):1–14, 2008.
- [Dav06] T. A. Davis. *Direct Methods for Sparse Linear Systems*. Society for Industrial and Applied Mathematics, 2006.
- [Dav09] T. A. Davis. The University of Florida Sparse Matrix Collection. Technical report, Dept. of Computer and Information Science and Engineering, Univ. of Florida, 2009.
- [dBE08] A. d'Aspremont, O. Banerjee, and L. El Ghaoui. First-order methods for sparse covariance selection. *SIAM Journal on Matrix Analysis and Applications*, 30(1):56–66, 2008.

- [DCHD90] J. J. Dongarra, J. Du Croz, S. Hammarling, and I. Duff. A set of level 3 basic linear algebra subprograms. *ACM Transactions on Mathematical Software*, 16(1):1–17, 1990.
- [DCHH88] J. J. Dongarra, J. Du Croz, S. Hammarling, and R. J. Hanson. An extended set of Fortran basic linear algebra subprograms. *ACM Transactions on Mathematical Software*, 14(1):1–17, 1988.
- [DR83] I. S. Duff and J. K. Reid. The multifrontal solution of indefinite sparse symmetric linear equations. *ACM Transactions on Mathematical Software*, 9(3):302–325, 1983.
- [DV09] J. Dahl and L. Vandenberghe. *CHOMPACT: Chordal Matrix Package*. abel.ee.ucla.edu/chompack, 2009.
- [DV10] J. Dahl and L. Vandenberghe. *CVXOPT: A Python Package for Convex Optimization*. abel.ee.ucla.edu/cvxopt, 2010.
- [DVR08] J. Dahl, L. Vandenberghe, and V. Roychowdhury. Covariance selection for non-chordal graphs via chordal embedding. *Optimization Methods and Software*, 23(4):501–520, 2008.
- [FHT08] J. Friedman, T. Hastie, and R. Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432, 2008.
- [FKMN00] M. Fukuda, M. Kojima, K. Murota, and K. Nakata. Exploiting sparsity in semidefinite programming via matrix completion I: general framework. *SIAM Journal on Optimization*, 11:647–674, 2000.
- [Fle95] R. Fletcher. An optimal positive definite update for sparse Hessian matrices. *SIAM Journal on Optimization*, 5(1):192–218, February 1995.
- [GJSW84] R. Grone, C. R. Johnson, E. M. Sá, and H. Wolkowicz. Positive definite completions of partial Hermitian matrices. *Linear Algebra and Appl.*, 58:109–124, 1984.
- [GL06] A. Guermouche and J.-Y. L’Excellent. Constructing memory-minimizing schedules for multifrontal methods. *ACM Transactions on Mathematical Software*, 32:17–32, March 2006.
- [GP80] G. H. Golub and R. J. Plemmons. Large-scale geodetic least-squares adjustment by dissection and orthogonal decomposition. *Linear Algebra and Its Applications*, 34(3):3–27, 1980.
- [GV96] G. H. Golub and C. F. Van Loan. *Matrix Computations*. Johns Hopkins University Press, 3rd edition, 1996.
- [GW08] A. Griewank and A. Walther. *Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation*. Society for Industrial and Applied Mathematics, second edition, 2008.
- [HLPL06] J. Z. Huang, N. Liu, M. Pourahmadi, and L. Liu. Covariance matrix selection and estimation via penalised normal likelihood. *Biometrika*, 93(1):85–98, 2006.

- [Lau01] M. Laurent. Matrix completion problems. In C. A. Floudas and P. M. Pardalos, editors, *Encyclopedia of Optimization*, volume III, pages 221–229. Kluwer, 2001.
- [Liu86] J. W. H. Liu. On the storage requirement in the out-of-core multifrontal method for sparse factorization. *ACM Transactions on Mathematical Software*, 12:249–264, 1986.
- [Liu88] J. W. H. Liu. Equivalent sparse matrix reordering by elimination tree rotations. *SIAM Journal on Scientific and Statistical Computing*, 9:424–444, May 1988.
- [Liu90] J. W. H. Liu. The role of elimination trees in sparse factorization. *SIAM Journal on Matrix Analysis and Applications*, 11:134–172, 1990.
- [Liu92] J. W. H. Liu. The multifrontal method for sparse matrix solution: theory and practice. *SIAM Review*, 34:82–109, 1992.
- [LNP93] J. W. H. Liu, E. G. Ng, and B. W. Peyton. On finding supernodes for sparse matrix computations. *SIAM Journal on Matrix Analysis and Applications*, 14(1):242–252, 1993.
- [LPP89] J. G. Lewis, B. W. Peyton, and A. Pothen. A fast algorithm for reordering sparse matrices for parallel factorization. *SIAM Journal on Scientific and Statistical Computing*, 10(6):1146–1173, 1989.
- [LT10] L. Li and K.-C. Toh. An inexact interior point method for L1-regularized sparse covariance selection. *Mathematical Programming Computation*, 2:291–315, 2010.
- [NFF⁺03] K. Nakata, K. Fujitsawa, M. Fukuda, M. Kojima, and K. Murota. Exploiting sparsity in semidefinite programming via matrix completion II: implementation and numerical details. *Mathematical Programming Series B*, 95:303–327, 2003.
- [PS90] A. Pothen and C. Sun. Compact clique tree data structures in sparse matrix factorizations. In T. F. Coleman and Y. Li, editors, *Large-Scale Numerical Optimization*, pages 180–204. Society for Industrial and Applied Mathematics, 1990.
- [SMG10] K. Scheinberg, S. Ma, and D. Goldfarb. Sparse inverse covariance selection via alternating linearization methods. In J. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R.S. Zemel, and A. Culotta, editors, *Advances in Neural Information Processing Systems 23*, pages 2101–2109. 2010.
- [SV04] G. Srijuntongsiri and S. A. Vavasis. A fully sparse implementation of a primal-dual interior-point potential reduction method for semidefinite programming, 2004.
- [Yam08] N. Yamashita. Sparse quasi-Newton updates with positive definite matrix completion. *Mathematical Programming, Series A*, 115(1):1–30, 2008.